

# Stochastic Differential Equations: Models and Numerics <sup>1</sup>

Jesper Carlsson      Kyoung-Sook Moon      Anders Szepessy  
Raúl Tempone      Georgios Zouraris

April 6, 2018

<sup>1</sup>This is a draft. Comments and improvements are welcome.

# Contents

<b>1</b>	<b>Introduction to Mathematical Models and their Analysis</b>	<b>4</b>
1.1	Noisy Evolution of Stock Values . . . . .	5
1.2	Molecular Dynamics . . . . .	6
1.3	Optimal Control of Investments . . . . .	7
1.4	Calibration of the Volatility . . . . .	8
1.5	The Coarse-graining and Discretization Analysis . . . . .	8
<b>2</b>	<b>Stochastic Integrals</b>	<b>11</b>
2.1	Probability Background . . . . .	11
2.2	Brownian Motion . . . . .	12
2.3	Approximation and Definition of Stochastic Integrals . . . . .	13
<b>3</b>	<b>Stochastic Differential Equations</b>	<b>23</b>
3.1	Approximation and Definition of SDE . . . . .	23
3.2	Itô's Formula . . . . .	30
3.3	Stratonovich Integrals . . . . .	35
3.4	Systems of SDE . . . . .	36
<b>4</b>	<b>The Feynman-Kâc Formula and the Black-Scholes Equation</b>	<b>38</b>
4.1	The Feynman-Kâc Formula . . . . .	38
4.2	Black-Scholes Equation . . . . .	39
<b>5</b>	<b>The Monte-Carlo Method</b>	<b>44</b>
5.1	Statistical Error . . . . .	44
5.2	Time Discretization Error . . . . .	48
<b>6</b>	<b>Finite Difference Methods</b>	<b>54</b>
6.1	American Options . . . . .	54
6.2	Lax Equivalence Theorem . . . . .	56
<b>7</b>	<b>The Finite Element Method and Lax-Milgram's Theorem</b>	<b>62</b>
7.1	The Finite Element Method . . . . .	63
7.2	Error Estimates and Adaptivity . . . . .	67
7.2.1	An A Priori Error Estimate . . . . .	67

7.2.2	An A Posteriori Error Estimate . . . . .	69
7.2.3	An Adaptive Algorithm . . . . .	70
7.3	Lax-Milgram's Theorem . . . . .	71
<b>8</b>	<b>Markov Chains, Duality and Dynamic Programming</b>	<b>76</b>
8.1	Introduction . . . . .	76
8.2	Markov Chains . . . . .	77
8.3	Expected Values . . . . .	79
8.4	Duality and Qualitative Properties . . . . .	81
8.5	Dynamic Programming . . . . .	83
8.6	Examples and Exercises . . . . .	85
<b>9</b>	<b>Optimal Control and Inverse Problems</b>	<b>87</b>
9.1	The Deterministic Optimal Control Setting . . . . .	88
9.1.1	Examples of Optimal Control . . . . .	89
9.1.2	Approximation of Optimal Control . . . . .	90
9.1.3	Motivation of the Lagrange formulation . . . . .	91
9.1.4	Dynamic Programming and the HJB Equation . . . . .	93
9.1.5	Characteristics and the Pontryagin Principle . . . . .	94
9.1.6	Generalized Viscosity Solutions of HJB Equations . . . . .	97
9.1.7	Maximum Norm Stability of Viscosity Solutions . . . . .	104
9.2	Numerical Approximation of ODE Constrained Minimization . . . . .	106
9.2.1	Optimization Examples . . . . .	108
9.2.2	Solution of the Discrete Problem . . . . .	117
9.2.3	Convergence of Euler Pontryagin Approximations . . . . .	120
9.2.4	How to obtain the Controls . . . . .	125
9.2.5	Inverse Problems and Tikhonov Regularization . . . . .	125
9.2.6	Smoothed Hamiltonian as a Tikhonov Regularization . . . . .	131
9.2.7	General Approximations . . . . .	133
9.3	Optimal Control of Stochastic Differential Equations . . . . .	135
9.3.1	An Optimal Portfolio . . . . .	136
9.3.2	Dynamic Programming and HJB Equations . . . . .	138
9.3.3	Relation of Hamilton-Jacobi Equations and Conservation Laws . . . . .	141
9.3.4	Numerical Approximations of Conservation Laws and Hamilton-Jacobi Equations . . . . .	144
<b>10</b>	<b>Rare Events and Reactions in SDE</b>	<b>148</b>
10.1	Invariant Measures and Ergodicity . . . . .	150
10.2	Reaction Rates . . . . .	155
10.3	Reaction Paths . . . . .	159

<b>11 Molecular dynamics</b>	<b>162</b>
11.1 Molecular dynamics at constant temperature: Zwanzig's model and derivation of Langevin dynamics . . . . .	163
11.2 The Gibbs distribution derived from dynamic stability . . . . .	165
11.3 Smoluchowski dynamics derived from Langevin dynamics . . . . .	168
11.4 Macroscopic conservation laws for compressible fluids motivated from molecular dynamics . . . . .	168
11.4.1 A general potential . . . . .	175
<b>12 Appendices</b>	<b>177</b>
12.1 Tomography Exercise . . . . .	177
12.2 Molecular Dynamics . . . . .	182
<b>13 Recommended Reading</b>	<b>194</b>

# Chapter 1

## Introduction to Mathematical Models and their Analysis

The goal of this course is to give useful understanding for solving problems formulated by stochastic differential equations models in science, engineering and mathematical finance. Typically, these problems require numerical methods to obtain a solution and therefore the course focuses on basic understanding of stochastic and partial differential equations to construct reliable and efficient computational methods.

Stochastic and deterministic differential equations are fundamental for the modeling in Science and Engineering. As the computational power increases, it becomes feasible to use more accurate differential equation models and solve more demanding problems: for instance to determine input data from fundamental principles, to optimally reconstruct input data using measurements or to find the optimal construction of a design. There are therefore two interesting computational sides of differential equations:

- the forward problem, to accurately determine solutions of differential equations for given data with minimal computational work and prescribed accuracy, and
- the inverse problem, to determine the input data for differential equations, from optimal estimates, based either on measurements or on computations with a more fundamental model.

The model can be stochastic by two reasons:

- if calibration of data implies this, as in financial mathematics, or
- if fundamental microscopic laws generate stochastic behavior when coarse-grained, as in molecular dynamics for chemistry, material science and biology.

An understanding of which model and method should be used in a particular situation requires some knowledge of both the model approximation error and the discretization error of the method. The optimal method clearly minimizes the computational work for given accuracy. Therefore it is valuable to know something about computational accuracy and work for different numerical models and methods, which lead us to error estimates

and convergence results. In particular, our study will take into account the amount of computational work for alternative mathematical models and numerical methods to solve a problem with a given accuracy.

## 1.1 Noisy Evolution of Stock Values

Let us consider a stock value denoted by the time dependent function  $S(t)$ . To begin our discussion, assume that  $S(t)$  satisfies the differential equation

$$\frac{dS}{dt} = a(t)S(t),$$

which has the solution

$$S(t) = e^{\int_0^t a(u)du} S(0).$$

Our aim is to introduce some kind of noise in the above simple model of the form  $a(t) = r(t) + \text{"noise"}$ , taking into account that we do not know precisely how the evolution will be. An example of a "noisy" model we shall consider is the stochastic differential equation

$$dS(t) = r(t)S(t)dt + \sigma S(t)dW(t), \quad (1.1)$$

where  $dW(t)$  will introduce noise in the evolution. To seek a solution for the above, the starting point will be the discretization

$$S_{n+1} - S_n = r_n S_n \Delta t_n + \sigma_n S_n \Delta W_n, \quad (1.2)$$

where  $\Delta W_n$  are independent normally distributed random variables with zero mean and variance  $\Delta t_n$ , i.e.  $E[\Delta W_n] = 0$  and  $Var[\Delta W_n] = \Delta t_n = t_{n+1} - t_n$ . As will be seen later on, equation (1.1) may have more than one possible interpretation, and the characterization of a solution will be intrinsically associated with the numerical discretization used to solve it.

We shall consider, among others, applications to option pricing problems. An European call option is a contract which gives the right, but not the obligation, to buy a stock for a fixed price  $K$  at a fixed future time  $T$ . The celebrated Black-Scholes model for the value  $f : (0, T) \times (0, \infty) \rightarrow \mathbb{R}$  of an option is the partial differential equation

$$\begin{aligned} \partial_t f + rs\partial_s f + \frac{\sigma^2 s^2}{2} \partial_s^2 f &= rf, \quad 0 < t < T, \\ f(s, T) &= \max(s - K, 0), \end{aligned} \quad (1.3)$$

where the constants  $r$  and  $\sigma$  denote the riskless interest rate and the volatility respectively. If the underlying stock value  $S$  is modeled by the stochastic differential equation (1.1) satisfying  $S(t) = s$ , the Feynmann-Kač formula gives the alternative probability representation of the option price

$$f(s, t) = E[e^{-r(T-t)} \max(S(T) - K, 0) | S(t) = s], \quad (1.4)$$

which connects the solution of a partial differential equation with the expected value of the solution of a stochastic differential equation. Although explicit exact solutions can be found in particular cases, our emphasis will be on general problems and numerical solutions. Those can arise from discretization of (1.3), by finite difference or finite elements methods, or from Monte Carlo methods based on statistical sampling of (1.4), with a discretization (1.2). Finite difference and finite element methods lead to a discrete system of equations substituting derivatives for difference quotients, e.g.

$$f_t \approx \frac{f(t_{n+1}) - f(t_n)}{\Delta t},$$

while the Monte Carlo method discretizes a probability space by substituting expected values with averages of finite samples, e.g.  $\{S(T, \omega_j)\}_{j=1}^M$  and

$$f(s, t) \approx \sum_{j=1}^M \frac{e^{-r(T-t)} \max(S(T, \omega_j) - K, 0)}{M}.$$

Which method is best? The solution depends on the problem to solve and we will carefully study qualitative properties of the numerical methods to understand the answer.

## 1.2 Molecular Dynamics

An example where the noise can be derived from fundamental principles is molecular dynamics, modeling e.g. reactions in chemistry and biology. Theoretically molecular systems can be modeled by the Schrödinger equation

$$i\partial_t \Psi = H\Psi$$

where the unknown  $\Psi$  is a wave function depending on time  $t$  and the variables of coordinates and spins of all,  $M$ , nuclei and,  $N$ , electrons in the problem; and  $H$  is the Hamiltonian precisely defined by well known fundamental constants of nature and the Coulomb interaction of all nuclei and electrons. An important issue is its high computational complexity for problems with more than a few nuclei, due to the high dimension of  $\Psi$  which is roughly in  $L^2(\mathbb{R}^{3(M+N)})$ , see [LB05]. Already simulation of a single water molecule requires a partial differential equation in 39 space dimensions, which is a demanding task to solve also with modern sparse approximation techniques.

A substantial dimensional reduction is obtained with Born-Oppenheimer approximation treating the nuclei as classical particles with the electrons in the ground state corresponding to the current nuclei positions. This approximation, derived from a WKB approximation for heavy nuclei mass (see Section 11), leads to *ab initio* molecular dynamics

$$\begin{aligned} \dot{x}_t &= v_t, \\ m\dot{v}_t &= -V'(x_t). \end{aligned} \tag{1.5}$$

To determine the nuclei dynamics and find the electron energy (input to  $V$ ) means now to solve a differential equation in  $\mathbb{R}^{6M}$  where at each time step the electron ground state

energy needs to be determined for the current nuclei configuration  $x^t$ , see [LB05, Fre02]. To simulate large systems with many particles requires some simplification of the expensive force calculation  $\partial_{x_i} V$  involving the current position  $x_t \in \mathbb{R}^{3M}$  of all nuclei.

The Hamiltonian system (1.5) is often further modified. For instance, equation (1.5) corresponds to simulate a problem with the number of particles, volume and total energy held constant. Simulation of a system with constant number of particles, volume and temperature are often done by using (1.5) and regularly rescaling the kinetic energy to meet the fixed temperature constraint, using so called thermostats. A mathematically attractive alternative to approximate a system in constant temperature is to solve the Langevin-Itô stochastic differential equation

$$\begin{aligned} dx_t &= v_t dt, \\ m dv_t &= -(V'(x_t) + \tau^{-1} v_t) dt + (2k_B T \tau^{-1})^{1/2} dW_t \end{aligned} \tag{1.6}$$

where  $T$  is the temperature,  $k_B$  the Boltzmann constant,  $W$  is a standard Wiener process in  $\mathbb{R}^{3M}$  and  $\tau$  is a relaxation time parameter (which can be determined from molecular dynamics simulation). The Langevin model (1.6) can be derived from the Schrödinger equation under certain assumptions, which is the subject of Sections ?? to ?. If diffusion is important in the problem under study, one would like to make long simulations on times of order at least  $\tau^{-1}$ . A useful observation to efficiently simulate longer time is the fact that for  $\tau \rightarrow 0+$  the solution  $x_{s/\tau}$  of the Langevin equation (??) converges to the solution  $\bar{x}_s$  solving the Smoluchowski equation, also called Brownian dynamics

$$d\bar{x}_s = -V'(\bar{x}_s) ds + (2k_B T)^{1/2} d\bar{W}_s, \tag{1.7}$$

set in the slower diffusion time scale  $s = \tau t$ . Here, for simplicity, the mass is assumed to be the same for all particles and normalized to  $m = 1$  and  $\bar{W}$  is again a standard Wiener process in  $\mathbb{R}^{3M}$ . The Smoluchowski model hence has the advantage to be able to approximate particle systems over longer time and reducing to half the problem dimension by eliminating the velocity variables. In Section 11.3 we analyze the weak approximation error  $x_{s/\tau} \rightarrow \bar{x}_s$ . The next step in the coarse-graining process is to derive partial differential equations – for the mass, momentum and energy of a continuum fluid – from Langevin or Smoluchowski molecular dynamics, which determines the otherwise unspecified pressure, viscosity and heat conductivity; Section ?? shows an example of such a coarse-graining process in the case of modelling a solid-liquid melt.

### 1.3 Optimal Control of Investments

Suppose that we invest in a risky asset, whose value  $S(t)$  evolves according to the stochastic differential equation  $dS(t) = \mu S(t) dt + \sigma S(t) dW(t)$ , and in a riskless asset  $Q(t)$  that evolves with  $dQ(t) = rQ(t) dt$ ,  $r < \mu$ . Our total wealth is then  $X(t) = Q(t) + S(t)$  and the goal is to determine an optimal instantaneous policy of investment in order to maximize the expected value of our wealth at a given final time  $T$ . Let the proportion of the total wealth invested on the risky asset at a given time  $t$ ,  $\alpha(t)$ , be defined by



$\alpha(t)X(t) = S(t)$ , so that  $(1 - \alpha(t))X(t) = Q(t)$  with  $\alpha(t) \in [0, 1]$ . Then our optimal control problem can be stated as

$$\max_{\alpha} E[g(X(T)) | X(t) = x] \equiv u(t, x),$$

where  $g$  is a given function. How can we determine an optimal  $\alpha$ ? The solution of this problem can be obtained by means of a Hamilton Jacobi equation, which is in general a nonlinear partial differential equation of the form

$$u_t + H(u, u_x, u_{xx}) = 0,$$

where  $H(u, u_x, u_{xx}) := \max_{\alpha} ((\mu\alpha x + r(1 - \alpha)x)u_x + \sigma^2\alpha^2x^2u_{xx}/2)$ . Part of our work is to study the theory of Hamilton Jacobi equations and numerical methods for control problems to determine the Hamiltonian  $H$  and the control  $\alpha$ . It turns out that typically the Hamiltonian needs to slightly modified in order to compute an approximate solution: Section 9 explains why and how. We call such modifications regularizations.

## 1.4 Calibration of the Volatility

Another important application of optimal control we will study is to solve inverse problems for differential equations in order to determine the input data for the differential equation from observed solution values, such as finding the volatility in the Black-Scholes equation from observed option prices: the option values can be used to determine the volatility function implicitly. The objective in the optimal control formulation is then to find a volatility function that yields option prices that deviate as little as possible from the measured option prices. The dynamics is the Black-Scholes equation with the volatility function to be determined, that is the dynamics is a deterministic partial differential equation and the volatility is the control function, see Section 9.2.1.1. This is a typical inverse problem: it is called inverse because in the standard view of the Black-Scholes equation relating the option values and the volatility, the option price is the unknown and the volatility is the data; while here the formulation is reversed with option prices as data and volatility as unknown in the same Black-Scholes equation. Inverse problems are often harder to solve than the forward problem and need to be regularized as explained in Section 9.

## 1.5 The Coarse-graining and Discretization Analysis

Our analysis of models and discretization methods use only one basic idea, which we present here for a deterministic problem of two differential equations

$$\dot{X}^t = a(X^t)$$

and

$$\dot{\bar{X}}^t = \bar{a}(\bar{X}^t).$$

We may think of the two given fluxes  $a$  and  $\bar{a}$  as either two different differential equation models or two discretization methods. The goal is to estimate a quantity of interest  $g(X^T)$ , e.g. the potential energy of a molecular dynamic system, the lift of an airfoil or the contract of a contingent claim in financial mathematics. Consider therefore a given function  $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$  with a solution  $X : [0, T] \rightarrow \mathbb{R}^d$ , e.g. the coordinates of atoms in a molecular system or a discretization of mass, momentum and energy of a fluid. To understand the global error  $g(X^T) - g(\bar{X}^T)$  we introduce the value function  $\bar{u}(x, t) := g(\bar{X}^T; \bar{X}^t = x)$ , which solves the partial differential equation

$$\begin{aligned} \partial_t \bar{u}(x, t) + \bar{a}(x) \partial_x \bar{u}(x, t) &= 0 \quad t < T \\ u(\cdot, T) &= g \end{aligned} \tag{1.8}$$

This definition and telescoping cancellation imply that the global error has the representation

$$\begin{aligned} g(X^T) - g(\bar{X}^T) &= \bar{u}(X^T, T) - \bar{u}(\underbrace{\bar{X}^0}_{=X^0}, 0) \\ &= \bar{u}(X^T, T) - \bar{u}(X^0, 0) \\ &= \int_0^T d\bar{u}(X^t, t) \\ &= \int_0^T \partial_t \bar{u}(X^t, t) + \dot{X}^t \partial_x \bar{u}(X^t, t) dt \\ &= \int_0^T \partial_t \bar{u}(X^t, t) + \bar{a}(X^t, t) \partial_x \bar{u}(X^t, t) dt \\ &= \int_0^T (-\bar{a}(X^t, t) + a(X^t, t)) \partial_x \bar{u}(X^t, t) dt. \end{aligned} \tag{1.9}$$

Here we can identify the local error in terms of the residual  $-\bar{a}(X^t, t) + a(X^t, t)$  multiplied by the weight  $\partial_x \bar{u}(X^t, t)$  and summed over all time steps. Note that the difference of the two solutions in the global error is converted into a weighted average of the residual  $-\bar{a}(X^t, t) + a(X^t, t)$  along only one solution  $X^t$ ; the representation is therefore the residual of  $X$ -path inserted into the  $\bar{u}$ -equation. We may view the error representation as a weak form of Lax Equivalence result, which states that the combination of consistence and stability imply convergence: consistence means that the flux  $\bar{a}$  approximates  $a$ ; stability means that  $\partial_x \bar{u}$  is bounded in some sense; and convergence means that the global error  $g(X^T) - g(\bar{X}^T)$  tends to zero. The equivalence, as it is usually known, is stated using bounds with appropriate norms and it has been the basis of the theoretical understanding of numerical methods.

The weak formulation (1.9) is easy to use and it is our basis for understanding both modelling and discretization errors. The weak form is particularly useful for estimating the weak approximation error, since it can take cancellation into account by considering the weaker concept of the value function instead of using absolute values and norms of differences of solution paths; the standard strong error analysis is obtained by estimating the norm of the difference of the two paths  $X$  and  $\bar{X}$ . Another attractive property of the weak representation (1.9) is that it can be applied both in *a priori* form to give

qualitative results, by combining it with analytical estimates of  $\partial_x \bar{u}$ , and in *a posteriori* form to obtain also quantitative results, by combining it with computer based estimates of  $\partial_x \bar{u}$ .

We first use the representation for understanding the weak approximation of stochastic differential equations and its time discretization, by extending the chain rule to Ito's formula and integrate over all outcomes (i.e. take the expected value). The value function solves a parabolic diffusion equation in this case, instead of the hyperbolic transport equation (1.8).

In the case of coarse-graining and modelling error, the representation is used for approximating

- Schrödinger dynamics by stochastic molecular Langevin dynamics,
- Kinetic Monte Carlo jump dynamics by SDE dynamics,
- Langevin dynamics by Smoluchowski dynamics, and
- Smoluchowski molecular dynamics by continuum phase-field dynamics.

We also use the representation for the important problem to analyse inverse problems, such as calibrating the volatility for stocks by observed option prices or finding an optimal portfolio of stocks and bonds. In an optimal control setting the extension is then to include a control parameter  $\alpha$  in the flux so that

$$\dot{X}^t = a(X^t, \alpha^t)$$

where the objective now is to find the minimum  $\min_{\alpha} g(X^T; X^t = x) =: u(x, t)$ . Then the value function  $u$  solves a nonlinear Hamilton-Jacobi-Bellman equation and the representation is extended by including a minimum over  $\alpha$ .

## Chapter 2

# Stochastic Integrals

This chapter introduces stochastic integrals, which will be the basis for stochastic differential equations in the next chapter. Here we construct approximations of stochastic integrals and prove an error estimate. The error estimate is then used to establish existence and uniqueness of stochastic integrals, which has the interesting ingredient of intrinsic dependence on the numerical approximation due to infinite variation. Let us first recall the basic definitions of probability we will use.

### 2.1 Probability Background

A probability space is a triple  $(\Omega, \mathcal{F}, P)$ , where  $\Omega$  is the set of outcomes,  $\mathcal{F}$  is the set of events and  $P : \mathcal{F} \rightarrow [0, 1]$  is a function that assigns probabilities to events satisfying the following definitions.

**Definition 2.1.** If  $\Omega$  is a given non empty set, then a  $\sigma$ -algebra  $\mathcal{F}$  on  $\Omega$  is a collection  $\mathcal{F}$  of subsets of  $\Omega$  that satisfy:

- (1)  $\Omega \in \mathcal{F}$ ;
- (2)  $F \in \mathcal{F} \Rightarrow F^c \in \mathcal{F}$ , where  $F^c = \Omega - F$  is the complement set of  $F$  in  $\Omega$ ; and
- (3)  $F_1, F_2, \dots \in \mathcal{F} \Rightarrow \bigcup_{i=1}^{+\infty} F_i \in \mathcal{F}$ .

**Definition 2.2.** A probability measure on  $(\Omega, \mathcal{F})$  is a set function  $P : \mathcal{F} \rightarrow [0, 1]$  such that:

- (1)  $P(\emptyset) = 0$ ,  $P(\Omega) = 1$ ; and
- (2) If  $A_1, A_2, \dots \in \mathcal{F}$  are mutually disjoint sets then

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

**Definition 2.3.** A random variable  $X$ , in the probability space  $(\Omega, \mathcal{F}, P)$ , is a function  $X : \Omega \rightarrow \mathbb{R}^d$  such that the inverse image

$$X^{-1}(A) \equiv \{\omega \in \Omega : X(\omega) \in A\} \in \mathcal{F},$$

for all open subsets  $A$  of  $\mathbb{R}^d$ .

**Definition 2.4** (Independence of random variables). Two sets  $A, B \in \mathcal{F}$  are said to be independent if

$$P(A \cap B) = P(A)P(B).$$

Two independent random variables  $X, Y$  in  $\mathbb{R}^d$  are independent if

$$X^{-1}(A) \text{ and } Y^{-1}(B) \text{ are independent for all open sets } A, B \subseteq \mathbb{R}^d.$$

**Definition 2.5.** A stochastic process  $X : [0, T] \times \Omega \rightarrow \mathbb{R}^d$  in the probability space  $(\Omega, \mathcal{F}, P)$  is a function such that  $X(t, \cdot)$  is a random variable in  $(\Omega, \mathcal{F}, P)$  for all  $t \in (0, T)$ . We will often write  $X(t) \equiv X(t, \cdot)$ .

The  $t$  variable will usually be associated with the notion of time.

**Definition 2.6.** Let  $X : \Omega \rightarrow \mathbb{R}$  be a random variable and suppose that the density function

$$p'(x) = \frac{P(X \in dx)}{dx}$$

is integrable. The expected value of  $X$  is then defined by the integral

$$E[X] = \int_{-\infty}^{\infty} xp'(x)dx, \tag{2.1}$$

which also can be written

$$E[X] = \int_{-\infty}^{\infty} xdp(x). \tag{2.2}$$

The last integral makes sense also in general when the density function is a measure, e.g. by successive approximation with random variables possessing integrable densities. A point mass, i.e. a Dirac delta measure, is an example of a measure.

**Exercise 2.7.** Show that if  $X, Y$  are independent random variables then

$$E[XY] = E[X]E[Y].$$

## 2.2 Brownian Motion

As a first example of a stochastic process, let us introduce

**Definition 2.8** (The Wiener process). The one-dimensional *Wiener process*  $W : [0, \infty) \times \Omega \rightarrow \mathbb{R}$ , also known as the Brownian motion, has the following properties:

- (1) with probability 1, the mapping  $t \mapsto W(t)$  is continuous and  $W(0) = 0$ ;  
(2) if  $0 = t_0 < t_1 < \dots < t_N = T$ , then the increments

$$W(t_N) - W(t_{N-1}), \dots, W(t_1) - W(t_0)$$

are *independent*; and

- (3) for all  $t > s$  the increment  $W(t) - W(s)$  has the *normal* distribution, with  $E[W(t) - W(s)] = 0$  and  $E[(W(t) - W(s))^2] = t - s$ , i.e.

$$P(W(t) - W(s) \in \Gamma) = \int_{\Gamma} \frac{e^{-\frac{y^2}{2(t-s)}}}{\sqrt{2\pi(t-s)}} dy, \quad \Gamma \subset \mathbb{R}.$$

Does there exist a Wiener process and how to construct  $W$  if it does? In computations we will only need to determine  $W$  at finitely many time steps  $\{t_n : n = 0, \dots, N\}$  of the form  $0 = t_0 < t_1 < \dots < t_N = T$ . The definition then shows how to generate  $W(t_n)$  by a sum of independent normal distributed random variables, see Example 2.20 for computational methods to generate independent normal distributed random variables. These independent increments will be used with the notation  $\Delta W_n = W(t_{n+1}) - W(t_n)$ . Observe, by Properties 1 and 3, that for fixed time  $t$  the Brownian motion  $W(t)$  is itself a normal distributed random variable. To generate  $W$  for all  $t \in \mathbb{R}$  is computationally infeasible, since it seems to require infinite computational work. Example 2.20 shows the existence of  $W$  by proving uniform convergence of successive continuous piecewise linear approximations. The approximations are based on an expansion in the orthogonal  $L^2(0, T)$  Haar-wavelet basis.

## 2.3 Approximation and Definition of Stochastic Integrals

**Remark 2.9** (Questions on the definition of a stochastic integral). Let us consider the problem of finding a reasonable definition for the stochastic integral  $\int_0^T W(t) dW(t)$ , where  $W(t)$  is the Wiener process. As a first step, let us discretize the integral by means of the *forward Euler* discretization

$$\sum_{n=0}^{N-1} W(t_n) \underbrace{(W(t_{n+1}) - W(t_n))}_{=\Delta W_n}.$$

Taking expected values we obtain by Property 2 of Definition 2.8

$$E\left[\sum_{n=0}^{N-1} W(t_n) \Delta W_n\right] = \sum_{n=0}^{N-1} E[W(t_n) \Delta W_n] = \sum_{n=0}^{N-1} E[W(t_n)] \underbrace{E[\Delta W_n]}_{=0} = 0.$$

Now let us use instead the *backward Euler* discretization

$$\sum_{n=0}^{N-1} W(t_{n+1}) \Delta W_n.$$

Taking expected values yields a different result:

$$\sum_{n=0}^{N-1} E[W(t_{n+1})\Delta W_n] = \sum_{n=0}^{N-1} E[W(t_n)\Delta W_n] + E[(\Delta W_n)^2] = \sum_{n=0}^{N-1} \Delta t = T \neq 0.$$

Moreover, if we use the *trapezoidal* method the result is

$$\begin{aligned} \sum_{n=0}^{N-1} E \left[ \frac{W(t_{n+1}) + W(t_n)}{2} \Delta W_n \right] &= \sum_{n=0}^{N-1} E[W(t_n)\Delta W_n] + E[(\Delta W_n)^2/2] \\ &= \sum_{n=0}^{N-1} \frac{\Delta t}{2} = T/2 \neq 0. \end{aligned}$$

□

Remark 2.9 shows that we need more information to define the stochastic integral  $\int_0^t W(s)dW(s)$  than to define a deterministic integral. We must decide if the solution we seek is the limit of the forward Euler method. In fact, limits of the forward Euler define the so called *Itô integral*, while the trapezoidal method yields the so called *Stratonovich integral*. It is useful to define the class of stochastic processes which can be Itô integrated. We shall restrict us to a class that allows computable quantities and gives convergence rates of numerical approximations. For simplicity, we begin with Lipschitz continuous functions in  $\mathbb{R}$  which satisfy (2.3) below. The next theorem shows that once the discretization method is fixed to be the forward Euler method, the discretizations converge in  $L^2$ . Therefore the limit of forward Euler discretizations is well defined, i.e. the limit does not depend on the sequence of time partitions, and consequently the limit can be used to define the Itô integral.

**Theorem 2.10.** *Suppose there exist a positive constant  $C$  such that  $f : [0, T] \times \mathbb{R} \rightarrow \mathbb{R}$  satisfies*

$$|f(t + \Delta t, W + \Delta W) - f(t, W)| \leq C(\Delta t + |\Delta W|). \quad (2.3)$$

*Consider two different partitions of the time interval  $[0, T]$*

$$\begin{aligned} \{\bar{t}_n\}_{n=0}^{\bar{N}}, \quad \bar{t}_0 = 0, \quad \bar{t}_{\bar{N}} = T, \\ \{\bar{\bar{t}}_m\}_{m=0}^{\bar{\bar{N}}}, \quad \bar{\bar{t}}_0 = 0, \quad \bar{\bar{t}}_{\bar{\bar{N}}} = T, \end{aligned}$$

*with the corresponding forward Euler approximations*

$$\bar{I} = \sum_{n=0}^{\bar{N}-1} f(\bar{t}_n, W(\bar{t}_n))(W(\bar{t}_{n+1}) - W(\bar{t}_n)), \quad (2.4)$$

$$\bar{\bar{I}} = \sum_{m=0}^{\bar{\bar{N}}-1} f(\bar{\bar{t}}_m, W(\bar{\bar{t}}_m))(W(\bar{\bar{t}}_{m+1}) - W(\bar{\bar{t}}_m)). \quad (2.5)$$

Let the maximum time step  $\Delta t_{max}$  be

$$\Delta t_{max} = \max \left[ \max_{0 \leq n \leq \bar{N}-1} \bar{t}_{n+1} - \bar{t}_n, \max_{0 \leq m \leq \bar{\bar{N}}-1} \bar{\bar{t}}_{m+1} - \bar{\bar{t}}_m \right].$$

Then

$$E[(\bar{I} - \bar{\bar{I}})^2] = \mathcal{O}(\Delta t_{max}). \quad (2.6)$$

*Proof.* It is useful to introduce the finer grid made of the union of the nodes on the two grids

$$\{t_k\} \equiv \{\bar{t}_n\} \cup \{\bar{\bar{t}}_m\}.$$

Then in that grid we can write

$$\bar{I} - \bar{\bar{I}} = \sum_k \Delta f_k \Delta W_k,$$

where  $\Delta f_k = f(\bar{t}_n, W(\bar{t}_n)) - f(\bar{\bar{t}}_m, W(\bar{\bar{t}}_m))$ ,  $\Delta W_k = W(t_{k+1}) - W(t_k)$  and the indices  $m, n$  satisfy  $t_k \in [\bar{t}_m, \bar{t}_{m+1})$  and  $t_k \in [\bar{t}_n, \bar{t}_{n+1})$ , as depicted in Figure 2.1.

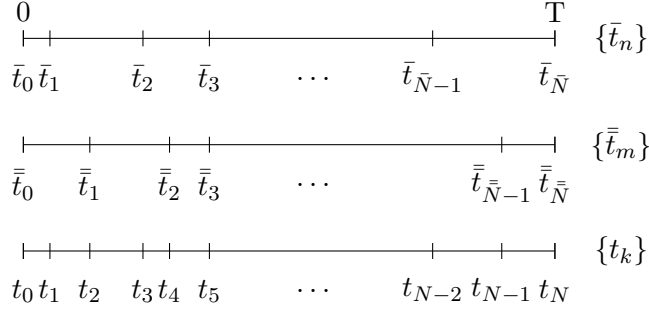


Figure 2.1: Mesh points used in the proof.

Therefore,

$$\begin{aligned} E[(\bar{I} - \bar{\bar{I}})^2] &= E\left[\sum_{k,l} \Delta f_k \Delta f_l \Delta W_l \Delta W_k\right] \\ &= 2 \sum_{k>l} \underbrace{E[\Delta f_k \Delta f_l \Delta W_l \Delta W_k]}_{=E[\Delta f_k \Delta f_l \Delta W_l]E[\Delta W_k]=0} + \sum_k E[(\Delta f_k)^2 (\Delta W_k)^2] \\ &= \sum_k E[(\Delta f_k)^2] E[(\Delta W_k)^2] = \sum_k E[(\Delta f_k)^2] \Delta t_k. \end{aligned} \quad (2.7)$$

Taking squares in (2.3) we arrive at  $|\Delta f_k|^2 \leq 2C^2((\Delta' t_k)^2 + (\Delta' W_k)^2)$  where  $\Delta' t_k = \bar{t}_n - \bar{\bar{t}}_m \leq \Delta t_{max}$  and  $\Delta' W_k = W(\bar{t}_n) - W(\bar{\bar{t}}_m)$ , using also the standard inequality



$(a + b)^2 \leq 2(a^2 + b^2)$ . Substituting this in (2.7) proves the theorem

$$\begin{aligned} E[(\bar{I} - \bar{\bar{I}})^2] &\leq \sum_k 2C^2 \left( (\Delta t_k)^2 + \underbrace{E[(\Delta W_k)^2]}_{=\Delta t_k} \right) \Delta t_k \\ &\leq 2C^2 T(\Delta t_{max}^2 + \Delta t_{max}). \end{aligned} \quad (2.8)$$

□

Thus, the sequence of approximations  $I_{\Delta t}$  is a Cauchy sequence in the Hilbert space of random variables generated by the norm  $\|I_{\Delta t}\|_{L^2} \equiv \sqrt{E[I_{\Delta t}^2]}$  and the scalar product  $(X, Y) \equiv E[XY]$ . The limit  $I$  of this Cauchy sequence defines the Itô integral

$$\sum_i f_i \Delta W_i \xrightarrow{L^2} I \equiv \int_0^T f(s, W(s)) dW(s).$$

**Remark 2.11** (Accuracy of strong convergence). If  $f(t, W(t)) = \bar{f}(t)$  is independent of  $W(t)$  we have first order convergence  $\sqrt{E[(\bar{I} - \bar{\bar{I}})^2]} = \mathcal{O}(\Delta t_{max})$ , whereas if  $f(t, W(t))$  depends on  $W(t)$  we only obtain one half order convergence  $\sqrt{E[(\bar{I} - \bar{\bar{I}})^2]} = \mathcal{O}(\sqrt{\Delta t_{max}})$ . The constant  $C$  in (2.3) and (2.9) measures the computational work to approximate the integral with the Euler method: to obtain an approximation error  $\epsilon$ , using uniform steps, requires by (2.8) the computational work corresponding to  $N = T/\Delta t \geq 4T^2C^2/\epsilon^2$  steps.

**Exercise 2.12.** Use the forward Euler discretization to show that

$$\int_0^T s dW(s) = TW(T) - \int_0^T W(s) ds$$

**Example 2.13** (Discrete Wiener process). A discrete Wiener process can be simulated by the following Octave/Matlab code:

```
% Simulation of Wiener process/Brownian path

N = 1E6;                % number of timesteps
randn('state',0);      % initialize random number generator
T = 1;                  % final time
dt = T/(N-1);          % time step
t = 0:dt:T;
dW = sqrt(dt)*randn(1,N-1); % Wiener increments
W = [0 cumsum(dW)];     % Brownian path
```

Brownian paths resulting from different seeds is shown in Figure 2.2, and in e.g. Exercise 2.12, the integrals can then be evaluated by

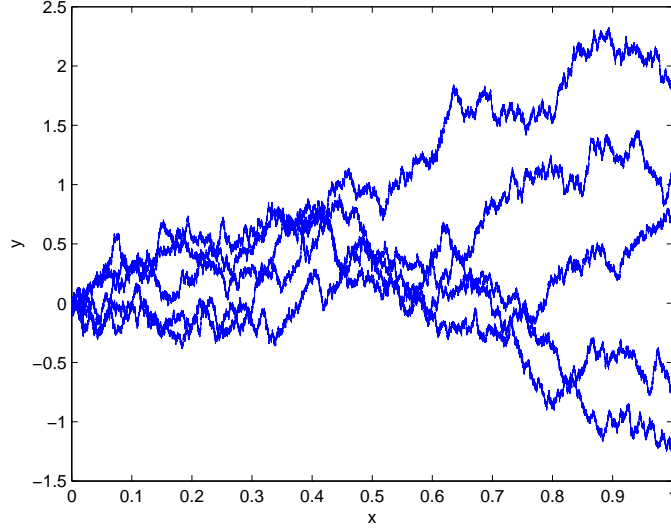


Figure 2.2: Brownian paths

```
LHS = sum(t(1:N-1).*dW);
RHS = T*W(N) - sum(W(1:N-1))*dt;
```

**Definition 2.14.** A process  $f : [0, T] \times \Omega \rightarrow \mathbb{R}$  is *adapted* if  $f(t, \cdot)$  only depends on events which are generated by  $W(s)$ ,  $s \leq t$ .

**Remark 2.15** (Extension to adapted Itô integration). Itô integrals can be extended to adapted processes. Assume  $f : [0, T] \times \Omega \rightarrow \mathbb{R}$  is adapted and that there is a constant  $C$  such that

$$\sqrt{E[|f(t + \Delta t, \omega) - f(t, \omega)|^2]} \leq C\sqrt{\Delta t}. \quad (2.9)$$

Then the proof of Theorem 2.10 shows that (2.4-2.6) still hold.

**Theorem 2.16** (Basic properties of Itô integrals).

Suppose that  $f, g : [0, T] \times \Omega \rightarrow \mathbb{R}$  are Itô integrable, e.g. adapted and satisfying (2.9), and that  $c_1, c_2$  are constants in  $\mathbb{R}$ . Then:

- (i)  $\int_0^T (c_1 f(s, \cdot) + c_2 g(s, \cdot)) dW(s) = c_1 \int_0^T f(s, \cdot) dW(s) + c_2 \int_0^T g(s, \cdot) dW(s)$ ,
- (ii)  $E \left[ \int_0^T f(s, \cdot) dW(s) \right] = 0$ ,
- (iii)  $E \left[ \left( \int_0^T f(s, \cdot) dW(s) \right) \left( \int_0^T g(s, \cdot) dW(s) \right) \right] = \int_0^T E [f(s, \cdot) g(s, \cdot)] ds$ .

*Proof.* To verify Property (ii), we first use that  $f$  is adapted and the independence of the increments  $\Delta W_n$  to show that for an Euler discretization

$$E\left[\sum_{n=0}^{N-1} f(t_n, \cdot) \Delta W_n\right] = \sum_{n=0}^{N-1} E[f(t_n, \cdot)] E[\Delta W_n] = 0.$$

It remains to verify that the limit of Euler discretizations preserves this property: Cauchy's inequality and the convergence result (2.6) imply that

$$\begin{aligned} |E\left[\int_0^T f(t, \cdot) dW(t)\right]| &= |E\left[\int_0^T f(t, \cdot) dW(t) - \sum_{n=0}^{N-1} f(t_n, \cdot) \Delta W_n\right] + \\ &\quad + E\left[\sum_{n=0}^{N-1} f(t_n, \cdot) \Delta W_n\right]| \\ &\leq \sqrt{E\left[\left(\int_0^T f(t, \cdot) dW(t) - \sum_{n=0}^{N-1} f(t_n, \cdot) \Delta W_n\right)^2\right]} \rightarrow 0. \end{aligned}$$

Property (i) and (iii) can be verified analogously.  $\square$

**Example 2.17** (The Monte-Carlo method). To verify Property (ii) in Theorem 2.16 numerically for some function  $f$  we can do a Monte-Carlo simulation where

$$\int_0^T f(s, \cdot) dW(s),$$

is calculated for several paths, or *realizations*, and then averaged:

```

% Monte-Carlo simulation

N = 1E3;                % number of timesteps
randn('state',0);      % initialize random number generator
T = 1;                 % final time
dt = T/N;              % time step
t = 0:dt:T;
M = 1E6;                % number of realisations
MC = zeros(1,M);       % vector to hold mean values

for i=1:M
    dW = sqrt(dt)*randn(1,N);    % Wiener increments
    W = [0 cumsum(dW)];          % Brownian paths
    f = t.^3.*sqrt(abs(W));      % some function
    int = sum(f(1:N).*dW);       % integral value
    if i==1
        MC(i) = int;
    else
        MC(i) = (MC(i-1)*(i-1)+int)/i; % new mean value
    end
end
end

```

In the above code the mean value of the integral is calculated for  $1, \dots, M$  realizations, and in Figure 2.3 we see that as the number of realizations grows, the mean value approaches zero as  $1/\sqrt{M}$ . Also, from the proof of Theorem 2.16 it can be seen that the number of time steps does not affect this convergence, so the provided code is inefficient, but merely serves as an illustration for the general case.

**Exercise 2.18.** Use the forward Euler discretization to show that

(a)  $\int_0^T W(s)dW(s) = \frac{1}{2}W(T)^2 - T/2.$

(b) Property (i) and (iii) in Theorem 2.16 hold.

**Exercise 2.19.** Consider the Ornstein-Uhlenbeck process defined by

$$X(t) = X_\infty + e^{-at}(X(0) - X_\infty) + b \int_0^t e^{-a(t-s)}dW(s), \quad (2.10)$$

where  $X_\infty, a$  and  $b$  are given real numbers. Use the properties of the Itô integral to compute  $E[X(t)], Var[X(t)], \lim_{t \rightarrow \infty} E[X(t)]$  and  $\lim_{t \rightarrow \infty} Var[X(t)]$ . Can you give an intuitive interpretation of the result?

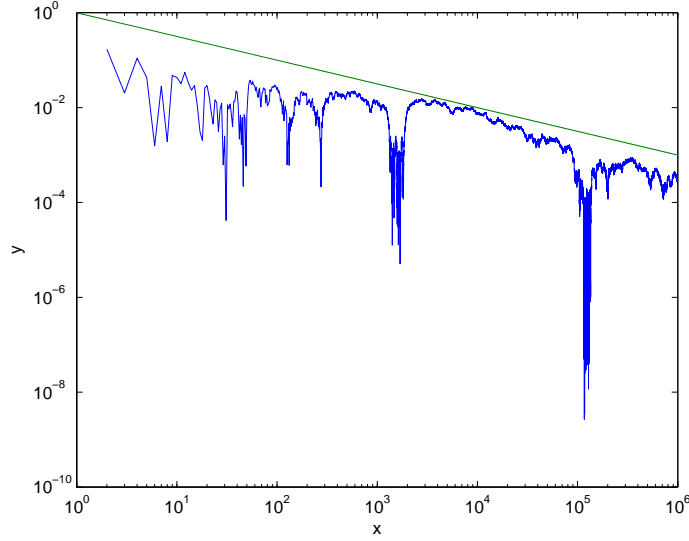


Figure 2.3: Absolute value of the mean for different number of realizations.

**Example 2.20** (Existence of a Wiener process). To construct a Wiener process on the time interval  $[0, T]$ , define the Haar-functions  $H_i$  by  $H_0(t) \equiv 1$  and for  $2^n \leq i < 2^{n+1}$  and  $n = 0, 1, 2, \dots$ , by

$$H_i(t) = \begin{cases} T^{-1/2}2^{n/2} & \text{if } (i - 2^n)2^{-n} \leq t/T < (i + 0.5 - 2^n)2^{-n}, \\ -T^{-1/2}2^{n/2} & \text{if } (i + 0.5 - 2^n)2^{-n} \leq t/T < (i + 1 - 2^n)2^{-n}, \\ 0 & \text{otherwise.} \end{cases} \quad (2.11)$$

Then  $\{H_i\}$  is an orthonormal basis of  $L^2(0, T)$ , (why?). Define the continuous piecewise linear function  $W^{(m)} : [0, T] \rightarrow \mathbb{R}$  by

$$W^{(m)}(t) = \sum_{i=1}^m \xi_i S_i(t), \quad (2.12)$$

where  $\xi_i$ ,  $i = 1, \dots, m$  are independent random variables with the normal distribution  $N(0, 1)$  and

$$S_i(t) = \int_0^t H_i(s) ds = \int_0^T 1_{(0,t)}(s) H_i(s) ds, \\ 1_{(0,t)}(s) = \begin{cases} 1 & \text{if } s \in (0, t), \\ 0 & \text{otherwise.} \end{cases}$$

The functions  $S_i$  are small "hat"-functions with a maximum value  $T^{-1/2}2^{-(n+2)/2}$  and zero outside an interval of length  $T2^{-n}$ . Let us postpone the proof that  $W^{(m)}$  converge uniformly and first assume this. Then the limit  $W(t) = \sum_{i=1}^{\infty} \xi_i S_i(t)$  is continuous. To

verify that the limit  $W$  is a Wiener process, we first observe that  $W(t)$  is a sum of normal distributed variables so that  $W(t)$  is also normal distributed. It remains to verify that the increments  $\Delta W_n$  and  $\Delta W_m$  are independent, for  $n \neq m$ , and  $E[(\Delta W_n)^2] = \Delta t_n$ . Parseval's equality shows the independence and the correct variance

$$\begin{aligned} E[\Delta W_n \Delta W_m] &= E\left[\sum_{i,j} \xi_i \xi_j (S_i(t_{n+1}) - S_i(t_n))(S_j(t_{m+1}) - S_j(t_m))\right] \\ &= \sum_{i,j} E[\xi_i \xi_j] (S_i(t_{n+1}) - S_i(t_n))(S_j(t_{m+1}) - S_j(t_m)) \\ &= \sum_i (S_i(t_{n+1}) - S_i(t_n))(S_i(t_{m+1}) - S_i(t_m)) \\ &\stackrel{\text{Parseval}}{=} \int_0^T 1_{(t_n, t_{n+1})}(s) 1_{(t_m, t_{m+1})}(s) ds = \begin{cases} 0 & \text{if } m \neq n, \\ t_{n+1} - t_n & \text{if } n = m. \end{cases} \end{aligned}$$

To prove uniform convergence, the goal is to establish

$$P\left(\sup_{t \in [0, T]} \sum_{i=1}^{\infty} |\xi_i| S_i(t) < \infty\right) = 1.$$

Fix a  $n$  and a  $t \in [0, T]$  then there is only one  $i$ , satisfying  $2^n \leq i < 2^{n+1}$ , such that  $S_i(t) \neq 0$ . Denote this  $i$  by  $i(t, n)$ . Let  $\chi_n \equiv \sup_{2^n \leq i < 2^{n+1}} |\xi_i|$ , then

$$\begin{aligned} \sup_{t \in [0, T]} \sum_{i=1}^{\infty} |\xi_i| S_i(t) &= \sup_{t \in [0, T]} \sum_{n=0}^{\infty} |\xi_{i(t, n)}| S_{i(t, n)}(t) \\ &\leq \sup_{t \in [0, T]} \sum_{n=0}^{\infty} |\xi_{i(t, n)}| T^{-1/2} 2^{-(n+2)/2} \\ &\leq \sum_{n=0}^{\infty} \chi_n T^{-1/2} 2^{-(n+2)/2}. \end{aligned}$$

If

$$\sum_{n=0}^{\infty} \chi_n 2^{-(n+2)/2} = \infty \tag{2.13}$$

on a set with positive probability, then  $\chi_n > n$  for infinitely many  $n$ , with positive probability, and consequently

$$\infty = E\left[\sum_{n=0}^{\infty} 1_{\{\chi_n > n\}}\right] = \sum_{n=0}^{\infty} P(\chi_n > n), \tag{2.14}$$

but

$$P(\chi_n > n) \leq P(\cup_{i=2^n}^{2^{n+1}} \{|\xi_i| > n\}) \leq 2^n P(|\xi_0| > n) \leq C 2^n e^{-n^2/4},$$

so that  $\sum_{n=0}^{\infty} P(\chi_n > n) < \infty$ , which contradicts (2.14) and (2.13). Therefore

$$P\left(\sup_{t \in [0, T]} \sum_{i=1}^{\infty} |\xi_i| S_i(t) < \infty\right) = 1,$$

which proves the uniform convergence.  $\square$

**Exercise 2.21** (Extension to multidimensional Itô integrals). The multidimensional Wiener process  $W$  in  $\mathbb{R}^l$  is defined by  $W(t) \equiv (W^1(t), \dots, W^l(t))$ , where  $W^i$ ,  $i = 1, \dots, l$  are independent one-dimensional Wiener processes. Show that

$$I_{\Delta t} \equiv \sum_{n=0}^{N-1} \sum_{i=1}^l f_i(t_n, \cdot) \Delta W_n^i$$

form a Cauchy sequence with  $E[(I_{\Delta t_1} - I_{\Delta t_2})^2] = \mathcal{O}(\Delta t_{max})$ , as in Theorem 2.10, provided  $f : [0, T] \times \Omega \rightarrow \mathbb{R}^l$  is adapted and (2.9) holds.

**Exercise 2.22.** Generalize Theorem 2.16 to multidimensional Itô integrals.

**Remark 2.23.** A larger class of Itô integrable functions are the functions in the Hilbert space

$$V = \left\{ f : [0, T] \times \Omega \rightarrow \mathbb{R}^l : f \text{ is adapted and } \int_0^T E[|f(t)|^2] dt < \infty \right\}$$

with the inner product  $\int_0^T E[f(t) \cdot g(t)] dt$ . This follows from the fact that every function in  $V$  can be approximated by adapted functions  $f_h$  that satisfy (2.9), for some constant  $C$  depending on  $h$ , so that  $\int_0^T E[|f(t, \cdot) - f_h(t, \cdot)|^2] dt \leq h$  as  $h \rightarrow 0$ . However, in contrast to Itô integration of the functions that satisfy (2.9), an approximation of the Itô integrals of  $f \in V$  does not in general give a convergence rate, but only convergence.

**Exercise 2.24.** Read Example 2.20 and show that the Haar-functions can be used to approximate stochastic integrals  $\int_0^T f(t) dW(t) \simeq \sum_{i=0}^m \xi_i f_i$ , for given deterministic functions  $f$  with  $f_i = \int_0^T f(s) H_i(s) ds$ . In what sense does  $dW(s) = \sum_{i=0}^{\infty} \xi_i H_i ds$  hold?

**Exercise 2.25.** Give an interpretation of the approximation (2.12) in terms of Brownian bridges, cf. [KS91].

## Chapter 3

# Stochastic Differential Equations

This chapter extends the work on stochastic integrals, in the last chapter, and constructs approximations of stochastic differential equations with an error estimate. Existence and uniqueness is then provided by the error estimate.

We will denote by  $C, C'$  positive constants, not necessarily the same at each occurrence.

### 3.1 Approximation and Definition of SDE

We will prove convergence of Forward Euler approximations of stochastic differential equations, following the convergence proof for Itô integrals. The proof is divided into four steps, including Grönwall's lemma below. The first step extends the Euler approximation  $\bar{X}(t)$  to all  $t \in [0, T]$ :

**Step 1.** Consider a grid in the interval  $[0, T]$  defined by the set of nodes  $\{\bar{t}_n\}_{n=0}^{\bar{N}}$ ,  $\bar{t}_0 = 0, \bar{t}_{\bar{N}} = T$  and define the discrete stochastic process  $\bar{X}$  by the forward Euler method

$$\bar{X}(\bar{t}_{n+1}) - \bar{X}(\bar{t}_n) = a(\bar{t}_n, \bar{X}(\bar{t}_n))(\bar{t}_{n+1} - \bar{t}_n) + b(\bar{t}_n, \bar{X}(\bar{t}_n))(W(\bar{t}_{n+1}) - W(\bar{t}_n)), \quad (3.1)$$

for  $n = 0, \dots, \bar{N} - 1$ . Now extend  $\bar{X}$  continuously, for theoretical purposes only, to all values of  $t$  by

$$\bar{X}(t) = \bar{X}(\bar{t}_n) + \int_{\bar{t}_n}^t a(\bar{t}_n, \bar{X}(\bar{t}_n))ds + \int_{\bar{t}_n}^t b(\bar{t}_n, \bar{X}(\bar{t}_n))dW(s), \quad \bar{t}_n \leq t < \bar{t}_{n+1}. \quad (3.2)$$

In other words, the process  $\bar{X} : [0, T] \times \Omega \rightarrow \mathbb{R}$  satisfies the stochastic differential equation

$$d\bar{X}(t) = \bar{a}(t, \bar{X})dt + \bar{b}(t, \bar{X})dW(t), \quad \bar{t}_n \leq t < \bar{t}_{n+1}, \quad (3.3)$$

where  $\bar{a}(t, \bar{X}) \equiv a(\bar{t}_n, \bar{X}(\bar{t}_n))$ ,  $\bar{b}(t, \bar{X}) \equiv b(\bar{t}_n, \bar{X}(\bar{t}_n))$ , for  $\bar{t}_n \leq t < \bar{t}_{n+1}$ , and the nodal values of the process  $\bar{X}$  is defined by the Euler method (3.1).

**Theorem 3.1.** *Let  $\bar{X}$  and  $\bar{X}$  be forward Euler approximations of the stochastic process  $X : [0, T] \times \Omega \rightarrow \mathbb{R}$ , satisfying the stochastic differential equation*

$$dX(t) = a(t, X(t))dt + b(t, X(t))dW(t), \quad 0 \leq t < T, \quad (3.4)$$



with time steps

$$\begin{aligned} \{\bar{t}_n\}_{n=0}^{\bar{N}}, \quad \bar{t}_0 = 0, \bar{t}_{\bar{N}} = T, \\ \{\bar{\bar{t}}_m\}_{m=0}^{\bar{\bar{N}}}, \quad \bar{\bar{t}}_0 = 0, \bar{\bar{t}}_{\bar{\bar{N}}} = T, \end{aligned}$$

respectively, and

$$\Delta t_{max} = \max \left[ \max_{0 \leq n \leq \bar{N}-1} \bar{t}_{n+1} - \bar{t}_n, \max_{0 \leq m \leq \bar{\bar{N}}-1} \bar{\bar{t}}_{m+1} - \bar{\bar{t}}_m \right].$$

Suppose that there exists a positive constant  $C$  such that the initial data and the given functions  $a, b : [0, T] \times \mathbb{R} \rightarrow \mathbb{R}$  satisfy

$$E[|\bar{X}(0)|^2 + |\bar{\bar{X}}(0)|^2] \leq C, \quad (3.5)$$

$$E\left[\left(\bar{X}(0) - \bar{\bar{X}}(0)\right)^2\right] \leq C\Delta t_{max}, \quad (3.6)$$

and

$$\begin{aligned} |a(t, x) - a(t, y)| &< C|x - y|, \\ |b(t, x) - b(t, y)| &< C|x - y|, \end{aligned} \quad (3.7)$$

$$|a(t, x) - a(s, x)| + |b(t, x) - b(s, x)| \leq C(1 + |x|)\sqrt{|t - s|}. \quad (3.8)$$

Then there is a constant  $K$  such that

$$\max \left\{ E[\bar{X}^2(t, \cdot)], E[\bar{\bar{X}}^2(t, \cdot)] \right\} \leq KT, \quad t < T, \quad (3.9)$$

and

$$E \left[ \left( \bar{X}(t, \cdot) - \bar{\bar{X}}(t, \cdot) \right)^2 \right] \leq K\Delta t_{max}, \quad t < T. \quad (3.10)$$

The basic idea for the extension of the convergence for Itô integrals to stochastic differential equations is

**Lemma 3.2** (Grönwall). *Assume that there exist positive constants  $A$  and  $K$  such that the function  $f : \mathbb{R} \rightarrow \mathbb{R}$  satisfies*

$$f(t) \leq K \int_0^t f(s) ds + A. \quad (3.11)$$

Then

$$f(t) \leq Ae^{Kt}.$$

*Proof.* Let  $I(t) \equiv \int_0^t f(s) ds$ . Then by (3.11)

$$\frac{dI}{dt} \leq KI + A,$$

and multiplying by  $e^{-Kt}$  we arrive at

$$\frac{d}{dt}(Ie^{-Kt}) \leq Ae^{-Kt}.$$

After integrating, and using  $I(0) = 0$ , we obtain  $I \leq A \frac{(e^{Kt}-1)}{K}$ . Substituting the last result in (3.11) concludes the proof.  $\square$

**Proof of the Theorem.** To prove (3.10), assume first that (3.9) holds. The proof is divided into the following steps:

- (1) Representation of  $\bar{X}$  as a process in continuous time: Step 1.
- (2) Use the assumptions (3.7) and (3.8).
- (3) Use the property (3) from Theorem 2.16.
- (4) Apply Grönwall's lemma.

**Step 2.** Consider another forward Euler discretization  $\bar{\bar{X}}$ , defined on a grid with nodes  $\{\bar{\bar{t}}_m\}_{m=0}^{\bar{\bar{N}}}$ , and subtract the two solutions to arrive at

$$\bar{X}(s) - \bar{\bar{X}}(s) \stackrel{(3.3)}{=} \bar{X}(0) - \bar{\bar{X}}(0) + \int_0^s \underbrace{(\bar{a} - \bar{\bar{a}})(t)}_{\equiv \Delta a(t)} dt + \int_0^s \underbrace{(\bar{b} - \bar{\bar{b}})(t)}_{\equiv \Delta b(t)} dW(t). \quad (3.12)$$

The definition of the discretized solutions implies that

$$\begin{aligned} \Delta a(t) &= (\bar{a} - \bar{\bar{a}})(t) = a(\bar{t}_n, \bar{X}(\bar{t}_n)) - a(\bar{\bar{t}}_m, \bar{\bar{X}}(\bar{\bar{t}}_m)) = \\ &= \underbrace{a(\bar{t}_n, \bar{X}(\bar{t}_n)) - a(t, \bar{X}(t))}_{=(I)} \\ &\quad + \underbrace{a(t, \bar{X}(t)) - a(t, \bar{\bar{X}}(t))}_{=(II)} \\ &\quad + \underbrace{a(t, \bar{\bar{X}}(t)) - a(\bar{\bar{t}}_m, \bar{\bar{X}}(\bar{\bar{t}}_m))}_{=(III)} \end{aligned}$$

where  $t \in [\bar{\bar{t}}_m, \bar{\bar{t}}_{m+1}) \cap [\bar{t}_n, \bar{t}_{n+1})$ , as shown in Figure 3.1. The assumptions (3.7) and (3.8) show that

$$\begin{aligned} |(I)| &\leq |a(\bar{t}_n, \bar{X}(\bar{t}_n)) - a(t, \bar{X}(\bar{t}_n))| + |a(t, \bar{X}(\bar{t}_n)) - a(t, \bar{X}(t))| \\ &\leq C|\bar{X}(\bar{t}_n) - \bar{X}(t)| + C(1 + |\bar{X}(\bar{t}_n)|)|t - \bar{t}_n|^{1/2}. \end{aligned} \quad (3.13)$$

Note that (3.7) and (3.8) imply

$$|a(t, x)| + |b(t, x)| \leq C(1 + |x|). \quad (3.14)$$

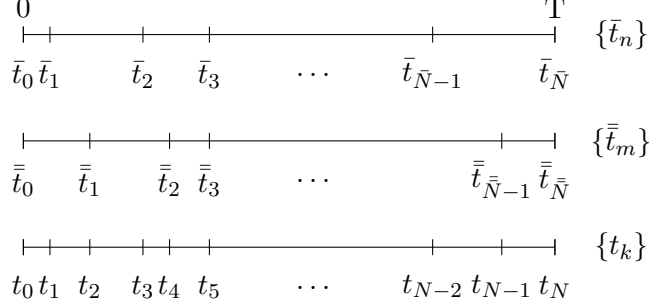


Figure 3.1: Mesh points used in the proof.

Therefore

$$\begin{aligned}
|\bar{X}(\bar{t}_n) - \bar{X}(t)| &\stackrel{(3.3)}{=} |a(\bar{t}_n, \bar{X}(\bar{t}_n))(t - \bar{t}_n) + b(\bar{t}_n, \bar{X}(\bar{t}_n))(W(t) - W(\bar{t}_n))| \\
&\stackrel{(3.14)}{\leq} C(1 + |\bar{X}(\bar{t}_n)|)((t - \bar{t}_n) + |W(t) - W(\bar{t}_n)|). \tag{3.15}
\end{aligned}$$

The combination of (3.13) and (3.15) shows

$$|(I)| \leq C(1 + |\bar{X}(\bar{t}_n)|) \left( |W(t) - W(\bar{t}_n)| + |t - \bar{t}_n|^{1/2} \right)$$

and in a similar way,

$$|(III)| \leq C(1 + |\bar{\bar{X}}(\bar{\bar{t}}_m)|) \left( |W(t) - W(\bar{\bar{t}}_m)| + |t - \bar{\bar{t}}_m|^{1/2} \right),$$

and by the assumptions (3.7)

$$|(II)| \stackrel{(3.7)}{\leq} C|\bar{X}(t) - \bar{\bar{X}}(t)|.$$

Therefore, the last three inequalities imply

$$\begin{aligned}
|\Delta a(t)|^2 &\leq (|(I)| + |(II)| + |(III)|)^2 \leq C_2 \left( |\bar{X}(t) - \bar{\bar{X}}(t)|^2 \right. \\
&\quad \left. + (1 + |\bar{X}(\bar{t}_n)|^2)(|t - \bar{t}_n| + |W(t) - W(\bar{t}_n)|^2) \right. \\
&\quad \left. + (1 + |\bar{\bar{X}}(\bar{\bar{t}}_m)|^2)(|t - \bar{\bar{t}}_m| + |W(t) - W(\bar{\bar{t}}_m)|^2) \right). \tag{3.16}
\end{aligned}$$

Recall that  $\max(t - \bar{t}_n, t - \bar{\bar{t}}_m) \leq \Delta t_{max}$ , and

$$E[(W(t) - W(s))^2] = t - s, \quad s < t,$$

so that the expected value of (3.16) and the assumption (3.9) yield

$$\begin{aligned}
E[|\Delta a(t)|^2] &\leq C \left( E[|\bar{X}(t) - \bar{\bar{X}}(t)|^2] + (1 + E[|\bar{X}(\bar{t}_n)|^2] + E[|\bar{\bar{X}}(\bar{t}_m)|^2]) \Delta t_{max} \right) \\
&\stackrel{(3.9)}{\leq} C \left( E[|\bar{X}(t) - \bar{\bar{X}}(t)|^2] + \Delta t_{max} \right). \tag{3.17}
\end{aligned}$$

Similarly, we have

$$E[|\Delta b(t)|^2] \leq C \left( E[|\bar{X}(t) - \bar{\bar{X}}(t)|^2] + \Delta t_{max} \right). \tag{3.18}$$

**Step 3.** Define a refined grid  $\{t_h\}_{h=0}^N$  by the union

$$\{t_h\} \equiv \{\bar{t}_n\} \cup \{\bar{\bar{t}}_m\}.$$

Observe that both the functions  $\Delta a(t)$  and  $\Delta b(t)$  are adapted and piecewise constant on the refined grid. The error representation (3.12) and (3) of Theorem 2.16 imply

$$\begin{aligned}
E[|\bar{X}(s) - \bar{\bar{X}}(s)|^2] &\leq E \left[ \left( \bar{X}(0) - \bar{\bar{X}}(0) + \int_0^s \Delta a(t) dt + \int_0^s \Delta b(t) dW(t) \right)^2 \right] \\
&\leq 3E[|\bar{X}(0) - \bar{\bar{X}}(0)|^2] \\
&\quad + 3E \left[ \left( \int_0^s \Delta a(t) dt \right)^2 \right] + 3E \left[ \left( \int_0^s \Delta b(t) dW(t) \right)^2 \right] \\
&\stackrel{(3.6)}{\leq} 3(C\Delta t_{max} + s \int_0^s E[(\Delta a(t))^2] dt + \int_0^s E[(\Delta b(t))^2] dt). \tag{3.19}
\end{aligned}$$

Inequalities (3.17-3.19) combine to

$$E[|\bar{X}(s) - \bar{\bar{X}}(s)|^2] \stackrel{(3.17-3.19)}{\leq} C \left( \int_0^s E[|\bar{X}(t) - \bar{\bar{X}}(t)|^2] dt + \Delta t_{max} \right). \tag{3.20}$$

**Step 4.** Finally, Grönwall's Lemma 3.2 applied to (3.20) implies

$$E[|\bar{X}(t) - \bar{\bar{X}}(t)|^2] \leq \Delta t_{max} C e^{Ct},$$

which finishes the proof. □

**Exercise 3.3.** Prove (3.9). Hint: Follow Steps 1-4 and use (3.5) .

**Corollary 3.4.** *The previous theorem yields a convergence result also in the  $L^2$  norm  $\|X\|^2 = \int_0^T E[X(t)^2] dt$ . The order of this convergence is  $1/2$ , i.e.  $\|\bar{X} - \bar{\bar{X}}\| = \mathcal{O}(\sqrt{\Delta t_{max}})$ .*

**Remark 3.5** (Strong and weak convergence). Depending on the application, our interest will be focused either on strong convergence

$$\|X(T) - \bar{X}(T)\|_{L^2[\Omega]} = \sqrt{E[(X(T) - \bar{X}(T))^2]} = \mathcal{O}(\sqrt{\Delta t}),$$

or on weak convergence  $E[g(X(T))] - E[g(\bar{X}(T))]$ , for given functions  $g$ . The next chapters will show first order convergence of expected values for the Euler method,

$$E[g(X(T)) - g(\bar{X}(T))] = \mathcal{O}(\Delta t),$$

and introduce Monte Carlo methods to approximate expected values  $E[g(\bar{X}(T))]$ . We will distinguish between strong and weak convergence by  $X_n \rightarrow X$ , denoting the strong convergence  $E[|X_n - X|^2] \rightarrow 0$  for random variables and  $\int_0^T E[|X_n(t) - X(t)|^2] dt \rightarrow 0$  for stochastic processes, and by  $X_n \rightharpoonup X$ , denoting the weak convergence  $E[g(X_n)] \rightarrow E[g(X)]$  for all bounded continuous functions  $g$ .

**Exercise 3.6.** Show that strong convergence,  $X_n \rightarrow X$ , implies weak convergence  $X_n \rightharpoonup X$ . Show also by an example that weak convergence,  $X_n \rightharpoonup X$ , does not imply strong convergence,  $X_n \rightarrow X$ . *Hint:* Let  $\{X_n\}$  be a sequence of independent identically distributed random variables.

Corollary 3.4 shows that successive refinements of the forward Euler approximation forms a Cauchy sequence in the Hilbert space  $V$ , defined by Definition 2.23. The limit  $X \in V$ , of this Cauchy sequence, satisfies the stochastic equation

$$X(s) = X(0) + \int_0^s a(t, X(t))dt + \int_0^s b(t, X(t))dW(t), \quad 0 < s \leq T, \quad (3.21)$$

and it is unique, (why?). Hence, we have constructed existence and uniqueness of solutions of (3.21) by forward Euler approximations. Let  $X$  be the solution of (3.21). From now on we use indistinctly also the notation

$$\begin{aligned} dX(t) &= a(t, X(t))dt + b(t, X(t))dW(t), \quad 0 < t \leq T \\ X(0) &= X_0. \end{aligned} \quad (3.22)$$

These notes focus on the Euler method to approximate stochastic differential equations (3.22). The following result motivates that there is no method with higher order convergence rate than the Euler method to control the strong error  $\int_0^1 E[(X(t) - \bar{X}(t))^2] dt$ , since even for the simplest equation  $dX = dW$  any linear approximation  $\hat{W}$  of  $W$ , based on  $N$  function evaluations, satisfies

**Theorem 3.7.** Let  $\hat{W}(t) = f(t, W(t_1), \dots, W(t_N))$  be any approximation of  $W(t)$ , which for fixed  $t$  is based on any linear function  $f(t, \cdot) : \mathbb{R}^N \rightarrow \mathbb{R}$ , and a partition  $0 = t_0 < \dots < t_N = 1$  of  $[0, 1]$ , then the strong approximation error is bounded from below by

$$\left( \int_0^1 E[(W(t) - \hat{W}(t))^2] dt \right)^{1/2} \geq \frac{1}{\sqrt{6N}}, \quad (3.23)$$

which is the same error as for the Euler method based on constant time steps and linear interpolation between the time steps.

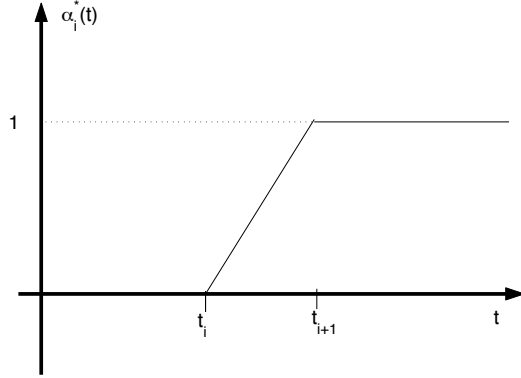


Figure 3.2: Optimal choice for weight functions  $\alpha_i$ .

*Proof.* The linearity of  $f(t, \cdot)$  implies that

$$\hat{W}(t) \equiv \sum_{i=1}^N \alpha_i(t) \Delta W_i$$

where  $\alpha_i : [0, 1] \rightarrow \mathbb{R}$ ,  $i = 1, \dots, N$  are any functions. The idea is to choose the functions  $\alpha_i : [0, 1] \rightarrow \mathbb{R}$ ,  $i = 1, \dots, N$  in an optimal way, and see that the minimum error satisfies (3.23). We have

$$\begin{aligned} & \int_0^1 E[(W(t) - \hat{W}(t))^2] dt \\ &= \int_0^1 (E[W^2(t)] - 2 \sum_{i=1}^N \alpha_i(t) E[W(t) \Delta W_i] + \sum_{i,j=1}^N \alpha_i(t) \alpha_j(t) E[\Delta W_i \Delta W_j]) dt \\ &= \int_0^1 t dt - 2 \int_0^1 \sum_{i=1}^N E[W(t) \Delta W_i] \alpha_i dt + \int_0^1 \sum_{i=1}^N \alpha_i^2(t) \Delta t_i dt \end{aligned}$$

and in addition

$$E[W(t) \Delta W_i] = \begin{cases} \Delta t_i, & t_{i+1} < t \\ (t - t_i), & t_i < t < t_{i+1} \\ 0, & t < t_i. \end{cases} \quad (3.24)$$

Perturbing the functions  $\alpha_i$ , to  $\alpha_i + \epsilon \delta_i$ ,  $\epsilon \ll 1$ , around the minimal value of  $\int_0^1 E[(W(t) - \hat{W}(t))^2] dt$  gives the following conditions for the optimum choice of  $\alpha_i$ , cf. Figure 3.2:

$$-2E[W(t) \Delta W_i] + 2\alpha_i^*(t) \Delta t_i = 0, \quad i = 1, \dots, N.$$

and hence

$$\begin{aligned}
\min \int_0^1 E[W(t) - \hat{W}(t)]^2 dt &= \int_0^1 t dt - \int_0^1 \sum_{i=1}^N \frac{E[W(t)\Delta W_i]^2}{\Delta t_i} dt \\
&\stackrel{(3.24)}{=} \sum_{n=1}^N (t_n + \Delta t_n/2)\Delta t_n - \sum_{n=1}^N \left( t_n \Delta t_n + \int_{t_n}^{t_{n+1}} \frac{(t - t_n)^2}{\Delta t_n} dt \right) \\
&= \sum_{n=1}^N (\Delta t_n)^2/6 \geq \frac{1}{6N}.
\end{aligned}$$

where Exercise 3.8 is used in the last inequality and proves the lower bound of the approximation error in the theorem. Finally, we note that by (3.24) the optimal  $\alpha_i^*(t) = \frac{E[W(t)\Delta W_i]}{\Delta t_i}$  is infact linear interpolation of the Euler method.  $\square$

**Exercise 3.8.** To verify the last inequality in the previous proof, compute

$$\begin{aligned}
&\min_{\Delta t} \sum_{n=1}^N (\Delta t_n)^2 \\
&\text{subject to} \\
&\sum_{n=1}^N (\Delta t_n) = 1.
\end{aligned}$$

## 3.2 Itô's Formula

Recall that using a forward Euler discretization we found the relation

$$\begin{aligned}
\int_0^T W(s)dW(s) &= W^2(T)/2 - T/2, \text{ or} \\
W(s)dW(s) &= d(W^2(s)/2) - ds/2,
\end{aligned} \tag{3.25}$$

whereas in the deterministic case we have  $y(s)dy(s) = d(y^2(s)/2)$ . The following useful theorem with Itô's formula generalizes (3.25) to general functions of solutions to the stochastic differential equations.

**Theorem 3.9.** *Suppose that the assumptions in Theorem 2.10 hold and that  $X$  satisfies the stochastic differential equation*

$$\begin{aligned}
dX(s) &= a(s, X(s))ds + b(s, X(s))dW(s), \quad s > 0 \\
X(0) &= X_0,
\end{aligned}$$

and let  $g : (0, +\infty) \times \mathbb{R} \rightarrow \mathbb{R}$  be a given bounded function in  $C^2((0, \infty) \times \mathbb{R})$ . Then  $y(t) \equiv g(t, X(t))$  satisfies the stochastic differential equation

$$\begin{aligned}
dy(t) &= \left( \partial_t g(t, X(t)) + a(t, X(t))\partial_x g(t, X(t)) + \frac{b^2(t, X(t))}{2}\partial_{xx}g(t, X(t)) \right) dt \\
&+ b(t, X(t))\partial_x g(t, X(t))dW(t),
\end{aligned} \tag{3.26}$$

*Proof.* We want to prove the Itô formula in the integral sense

$$\begin{aligned} & g(\tau, X(\tau)) - g(0, X(0)) \\ &= \int_0^\tau \left( \partial_t g(t, X(t)) + a(s, X(s)) \partial_x g(t, X(t)) + \frac{b^2(t, X(t))}{2} \partial_{xx} g(t, X(t)) \right) dt \\ & \quad + \int_0^\tau b(t, X(t)) \partial_x g(t, X(t)) dW(t). \end{aligned}$$

Let  $\bar{X}$  be a forward Euler approximation (3.1) and (3.2) of  $X$ , so that

$$\Delta \bar{X} \equiv \bar{X}(t_n + \Delta t_n) - \bar{X}(t_n) = a(t_n, \bar{X}(t_n)) \Delta t_n + b(t_n, \bar{X}(t_n)) \Delta W_n. \quad (3.27)$$

Taylor expansion of  $g$  up to second order gives

$$\begin{aligned} & g(t_n + \Delta t_n, \bar{X}(t_n + \Delta t_n)) - g(t_n, \bar{X}(t_n)) \\ &= \partial_t g(t_n, \bar{X}(t_n)) \Delta t_n + \partial_x g(t_n, \bar{X}(t_n)) \Delta \bar{X}(t_n) \\ & \quad + \frac{1}{2} \partial_{tt} g(t_n, \bar{X}(t_n)) \Delta t_n^2 + \partial_{tx} g(t_n, \bar{X}(t_n)) \Delta t_n \Delta \bar{X}(t_n) \\ & \quad + \frac{1}{2} \partial_{xx} g(t_n, \bar{X}(t_n)) (\Delta \bar{X}(t_n))^2 + o(\Delta t_n^2 + |\Delta \bar{X}_n|^2). \end{aligned} \quad (3.28)$$

The combination of (3.27) and (3.28) shows

$$\begin{aligned} & g(t_m, \bar{X}(t_m)) - g(0, \bar{X}(0)) = \sum_{n=0}^{m-1} (g(t_n + \Delta t_n, \bar{X}(t_n + \Delta t_n)) - g(t_n, \bar{X}(t_n))) \\ &= \sum_{n=0}^{m-1} \partial_t g \Delta t_n + \sum_{n=0}^{m-1} (\bar{a} \partial_x g \Delta t_n + \bar{b} \partial_x g \Delta W_n) + \frac{1}{2} \sum_{n=0}^{m-1} (\bar{b})^2 \partial_{xx} g (\Delta W_n)^2 \\ & \quad + \sum_{n=0}^{m-1} \left( (\bar{b} \partial_{tx} g + \bar{a} \bar{b} \partial_{xx} g) \Delta t_n \Delta W_n + \left( \frac{1}{2} \partial_{tt} g + \bar{a} \partial_{tx} g + \frac{1}{2} \bar{a}^2 \partial_{xx} g \right) \Delta t_n^2 \right) \\ & \quad + \sum_{n=0}^{m-1} o(\Delta t_n^2 + |\Delta \bar{X}(t_n)|^2). \end{aligned} \quad (3.29)$$

Let us first show that

$$\sum_{n=0}^{m-1} \bar{b}^2 \partial_{xx} g(\bar{X})(\Delta W_n)^2 \rightarrow \int_0^t b^2 \partial_{xx} g(X) ds,$$

as  $\Delta t_{max} \rightarrow 0$ . It is sufficient to establish

$$Y \equiv \frac{1}{2} \sum_{n=0}^{m-1} (\bar{b})^2 \partial_{xx} g((\Delta W_n)^2 - \Delta t_n) \rightarrow 0, \quad (3.30)$$

since (3.10) implies  $\sum_{n=0}^{m-1} (\bar{b})^2 \partial_{xx} g \Delta t_n \rightarrow \int_0^t b^2 \partial_{xx} g ds$ . Use the notation  $\alpha_i = ((\bar{b})^2 \partial_{xx} g)(t_i, \bar{X}(t_i))$  and independence to obtain



$$\begin{aligned}
E[Y^2] &= \sum_{i,j} E[\alpha_i \alpha_j ((\Delta W_i)^2 - \Delta t_i)((\Delta W_j)^2 - \Delta t_j)] \\
&= 2 \sum_{i>j} E[\alpha_i \alpha_j ((\Delta W_j)^2 - \Delta t_j)((\Delta W_i)^2 - \Delta t_i)] + \sum_i E[\alpha_i^2 ((\Delta W_i)^2 - \Delta t_i)^2] \\
&= 2 \sum_{i>j} E[\alpha_i \alpha_j ((\Delta W_j)^2 - \Delta t_j)] \underbrace{E[((\Delta W_i)^2 - \Delta t_i)]}_{=0} \\
&\quad + \sum_i E[\alpha_i^2] \underbrace{E[((\Delta W_i)^2 - \Delta t_i)^2]}_{=2\Delta t_i^2} \rightarrow 0,
\end{aligned}$$

when  $\Delta t_{max} \rightarrow 0$ , therefore (3.30) holds. Similar analysis with the other terms in (3.29) concludes the proof.  $\square$

**Remark 3.10.** The preceding result can be remembered intuitively by a Taylor expansion of  $g$  up to second order

$$dg = \partial_t g dt + \partial_x g dX + \frac{1}{2} \partial_{xx} g (dX)^2$$

and the relations:  $dt dt = dt dW = dW dt = 0$  and  $dW dW = dt$ .

**Example 3.11.** Let  $X(t) = W(t)$  and  $g(x) = \frac{x^2}{2}$ . Then

$$d\left(\frac{W^2(s)}{2}\right) = W(s)dW(s) + 1/2(dW(s))^2 = W(s)dW(s) + ds/2.$$

**Exercise 3.12.** Let  $X(t) = W(t)$  and  $g(x) = x^4$ . Verify that

$$d(W^4(s)) = 6W^2(s)ds + 4W^3(s)dW(s)$$

and

$$\frac{d}{ds}(E[g(W(s))]) = \frac{d}{ds}(E[(W(s))^4]) = 6s.$$

Apply the last result to compute  $E[W^4(t)]$  and  $E[(W^2(t) - t)^2]$ .

**Exercise 3.13.** Generalize the previous exercise to determine  $E[W^{2n}(t)]$ .

**Example 3.14.** We want to compute  $\int_0^T t dW(t)$ . Take  $g(t, x) = tx$ , and again  $X(t) = W(t)$ , so that

$$tW(t) = \int_0^t s dW(s) + \int_0^t W(s) ds$$

and finally  $\int_0^t s dW(s) = tW(t) - \int_0^t W(s) ds$ .

**Exercise 3.15.** Consider the stochastic differential equation

$$dX(t) = -a(X(t) - X_\infty)dt + b dW(t),$$

with initial data  $X(0) = X_0 \in \mathbb{R}$  and given  $a, b \in \mathbb{R}$ .

(i) Using that

$$X(t) - X(0) = -a \int_0^t (X(s) - X_\infty) ds + bW(t),$$

take the expected value and find an ordinary differential equation for the function  $m(t) \equiv E[X(t)]$ .

(ii) Use Itô's formula to find the differential of  $(X(t))^2$  and apply similar ideas as in (i) to compute  $Var[X(t)]$ .

(iii) Use an integrating factor to derive the exact solution (2.10) in Example 2.19. Compare your results from (i) and (ii) with this exact solution.

**Example 3.16.** Consider the stochastic differential equation

$$dS(t) = rS(t)dt + \sigma S(t)dW(t),$$

used to model the evolution of stock values. The values of  $r$  (interest rate) and  $\sigma$  (volatility) are assumed to be constant. Our objective is to find a closed expression for the solution, often called *geometric Brownian motion*. Let  $g(x) = \ln(x)$ . Then a direct application of Itô formula shows

$$d \ln(S(t)) = dS(t)/S(t) - 1/2 \left( \frac{\sigma^2 S^2(t)}{S^2(t)} \right) dt = rdt - \frac{\sigma^2}{2} dt + \sigma dW(t),$$

so that

$$\ln \left( \frac{S(T)}{S(0)} \right) = rT - \frac{T\sigma^2}{2} + \sigma W(T)$$

and consequently

$$S(T) = e^{(r - \frac{\sigma^2}{2})T + \sigma W(T)} S(0). \quad (3.31)$$

**Example 3.17** (Verification of strong and weak convergence). From the explicit formula (3.31) we can numerically verify the results on strong and weak convergence, given in Remark 3.5 for the Euler method. In the following code we calculate the strong and weak error by comparing the Euler simulation and the explicit value (3.31) at final time for several realizations. This is then tested for different time steps and the result in Figure 3.3 confirms a strong convergence of order 1/2 and a weak convergence of order 1.

```

% Strong and weak convergence for the Euler method

steps = [1:6];
for i=steps
    N = 2^i % number of timesteps
    randn('state',0);
    T = 1; dt = T/N; t = 0:dt:T;
    r = 0.1; sigma = 0.5; S0 = 100;
    M = 1E6; % number of realisations
    S = S0*ones(M,1); % S(0) for all realizations
    W = zeros(M,1); % W(0) for all realizations
    for j=1:N
        dW = sqrt(dt)*randn(M,1); % Wiener increments
        S = S + S.*(r*dt+sigma*dW); % processes at next time step
        W = W + dW; % Brownian paths at next step
    end
    ST = S0*exp( (r-sigma^2/2)*T + sigma*W ); % exact final value
    wError(i) = mean(S-ST); % weak error
    sError(i) = sqrt(mean((S-ST).^2)); % strong error
end
dt = T./2^steps;
loglog(dt,abs(wError),'o--',dt,dt,'--',dt,abs(sError),'o-',dt,sqrt(dt))

```

**Exercise 3.18.** Suppose that we want to simulate  $S(t)$ , defined in the previous example by means of the forward Euler method, i.e.

$$S_{n+1} = (1 + r\Delta t_n + \sigma\Delta W_n)S_n, \quad n = 0, \dots, N$$

As with the exact solution  $S(t)$ , we would like to have  $S_n$  positive. Then we could choose the time step  $\Delta t_n$  to reduce the probability of hitting zero

$$P(S_{n+1} < 0 | S_n = s) < \epsilon \ll 1. \quad (3.32)$$

Motivate a choice for  $\epsilon$  and find then the largest  $\Delta t_n$  satisfying (3.32).

**Remark 3.19.** The Wiener process has unbounded variation i.e.

$$E \left[ \int_0^T |dW(s)| \right] = +\infty.$$

This is the reason why the forward and backward Euler methods give different results.

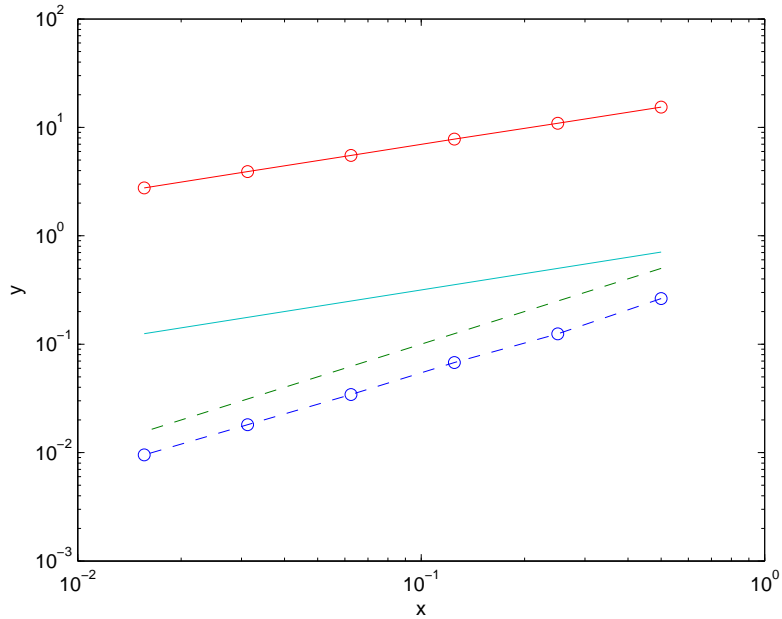


Figure 3.3: Strong and weak convergence.

We have for a uniform mesh  $\Delta t = T/N$

$$\begin{aligned}
 E\left[\sum_{i=0}^{N-1} |\Delta W_i|\right] &= \sum_{i=0}^{N-1} E[|\Delta W_i|] = \sum_{i=0}^{N-1} \sqrt{\frac{2\Delta t_i}{\pi}} \\
 &= \sqrt{\frac{2T}{\pi}} \sum_{i=0}^{N-1} \sqrt{1/N} = \sqrt{\frac{2NT}{\pi}} \rightarrow \infty, \quad \text{as } N \rightarrow \infty.
 \end{aligned}$$

### 3.3 Stratonovich Integrals

Recall from Chapter 2 that Itô integrals are constructed via forward Euler discretizations and Stratonovich integrals via the trapezoidal method, see Exercise 3.20. Our goal here is to express a Stratonovich integral

$$\int_0^T g(t, X(t)) \circ dW(t)$$

in terms of an Itô integral. Assume then that  $X(t)$  satisfies the Itô differential equation

$$dX(t) = a(t, X(t))dt + b(t, X(t))dW(t).$$

Then the relation reads

$$\begin{aligned} \int_0^T g(t, X(t)) \circ dW(t) &= \int_0^T g(t, X(t)) dW(t) \\ &+ \frac{1}{2} \int_0^T \partial_x g(t, X(t)) b(t, X(t)) dt. \end{aligned} \quad (3.33)$$

Therefore, Stratonovich integrals satisfy

$$dg(t, X(t)) = \partial_t g(t, X(t)) dt + \partial_x g(t, X(t)) \circ dX(t), \quad (3.34)$$

just like in the usual calculus.

**Exercise 3.20.** Use that Stratonovich integrals  $g(t, X(t)) \circ dW(t)$  are defined by limits of the trapezoidal method to verify (3.33), cf. Remark 2.9.

**Exercise 3.21.** Verify the relation (3.34), and use this to show that  $dS(t) = rS(t)dt + \sigma S(t) \circ dW(t)$  implies  $S(t) = e^{rt + \sigma W(t)} S(0)$ .

**Remark 3.22** (Stratonovich as limit of piecewise linear interpolations). Let  $R^N(t) \equiv W(t_n) + \frac{W(t_{n+1}) - W(t_n)}{t_{n+1} - t_n} (t - t_n)$ ,  $t \in (t_n, t_{n+1})$  be a piecewise linear interpolation of  $W$  on a given grid, and define  $X^N$  by  $dX^N(t) = a(X^N(t))dt + b(X^N(t))dR^N(t)$ . Then  $X^N \rightarrow X$  in  $L^2$ , where  $X$  is the solution of the Stratonovich stochastic differential equation

$$dX(t) = a(X(t))dt + b(X(t)) \circ dW(t).$$

In the special case when  $a(x) = rx$  and  $b(x) = \sigma x$  this follows from

$$d(\ln(X^N(t))) = rdt + \sigma dR^N,$$

so that

$$X^N(t) = e^{rt + \sigma R^N(t)} X(0).$$

The limit  $N \rightarrow \infty$  implies  $X^N(t) \rightarrow X(t) = e^{rt + \sigma W(t)} X(0)$ , as in Exercise 3.21.

### 3.4 Systems of SDE

Let  $W_1, W_2, \dots, W_l$  be scalar independent Wiener processes. Consider the  $l$ -dimensional Wiener process  $W = (W_1, W_2, \dots, W_l)$  and  $X : [0, T] \times \Omega \rightarrow \mathbb{R}^d$  satisfying for given drift  $a : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  and diffusion  $b : [0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^{d \times l}$  the Itô stochastic differential equation

$$dX_i(t) = a_i(t, X(t))dt + b_{ij}(t, X(t))dW_j(t), \text{ for } i = 1 \dots d. \quad (3.35)$$

Here and below we use of the summation convention

$$\alpha_j \beta_j \equiv \sum_j \alpha_j \beta_j,$$

i.e., if the same summation index appears twice in a term, the term denotes the sum over the range of this index. Theorem 3.9 can be directly generalized to the system (3.35).

**Theorem 3.23** (Itô 's formula for systems). *Let*

$$dX_i(t) = a_i(t, X(t))dt + b_{ij}(t, X(t))dW_j(t), \text{ for } i = 1 \dots d,$$

*and consider a smooth and bounded function  $g : \mathbb{R}^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$ . Then*

$$\begin{aligned} dg(t, X(t)) = & \left\{ \partial_t g(t, X(t)) + \partial_{x_i} g(t, X(t)) a_i(t, X(t)) \right. \\ & \left. + \frac{1}{2} b_{ik}(t, X(t)) \partial_{x_i x_j} g(t, X(t)) b_{jk}(t, X(t)) \right\} dt \\ & + \partial_{x_i} g(t, X(t)) b_{ij}(t, X(t)) dW_j(t), \end{aligned}$$

*or in matrix vector notation*

$$\begin{aligned} dg(t, X(t)) = & \left\{ \partial_t g(t, X(t)) + \nabla_x g(t, X(t)) a(t, X(t)) \right. \\ & \left. + \frac{1}{2} \text{trace} (b(t, X(t)) b^T(t, X(t)) \nabla_x^2 g(t, X(t))) \right\} dt \\ & + \nabla_x g(t, X(t)) b(t, X(t)) dW(t). \end{aligned}$$

**Remark 3.24.** The formal rules to remember Theorem 3.23 are Taylor expansion to second order and

$$\begin{aligned} dW_j dt &= dt dt = 0 \\ dW_i dW_j &= \delta_{ij} dt = \begin{cases} dt & \text{if } i = j, \\ 0 & \text{otherwise.} \end{cases} \end{aligned} \tag{3.36}$$

**Exercise 3.25.** Verify Remark 3.24.

## Chapter 4

# The Feynman-K ac Formula and the Black-Scholes Equation

### 4.1 The Feynman-K ac Formula

**Theorem 4.1.** *Suppose that  $a, b$  and  $g$  are smooth and bounded functions. Let  $X$  be the solution of the stochastic differential equation,*

$$dX(t) = a(t, X(t))dt + b(t, X(t))dW(t),$$

and let  $u(x, t) = E[g(X(T)) | X(t) = x]$ . Then  $u$  is the solution of the Kolmogorov backward equation

$$\begin{aligned} L^*u &\equiv u_t + au_x + \frac{1}{2}b^2u_{xx} = 0, \quad t < T \\ u(x, T) &= g(x). \end{aligned} \tag{4.1}$$

*Proof.* Define  $\hat{u}$  to be the solution of (4.1), i.e.  $L^*\hat{u} = 0$ ,  $\hat{u}(\cdot, T) = g(\cdot)$ . We want to verify that  $\hat{u}$  is the expected value  $E[g(X(T)) | X(t) = x]$ . The Itô formula applied to  $\hat{u}(X(t), t)$  shows

$$\begin{aligned} d\hat{u}(X(t), t) &= \left( \hat{u}_t + a\hat{u}_x + \frac{1}{2}b^2\hat{u}_{xx} \right) dt + b\hat{u}_x dW \\ &= L^*\hat{u}dt + b\hat{u}_x dW. \end{aligned}$$

Integrate this from  $t$  to  $T$  and use  $L^*\hat{u} = 0$  to obtain

$$\begin{aligned} \hat{u}(X(T), T) - \hat{u}(X(t), t) &= g(X(T)) - \hat{u}(X(t), t) \\ &= \int_t^T b\hat{u}_x dW(s). \end{aligned}$$

Take the expectation and use that the expected value of the Itô integral is zero,

$$\begin{aligned} E[g(X(T)) | X(t) = x] - \hat{u}(x, t) &= E\left[\int_t^T b(s, X(s))\hat{u}_x(X(s), s)dW(s) | X(t) = x\right] \\ &= 0. \end{aligned}$$

Therefore

$$\hat{u}(x, t) = E[g(X(T)) | X(t) = x],$$

which proves the theorem since the solution of Equation (4.1) is unique.  $\square$

**Exercise 4.2** (Maximum Principle). Let the function  $u$  satisfy

$$\begin{aligned} u_t + au_x + \frac{1}{2}b^2u_{xx} &= 0, \quad t < T \\ u(x, T) &= g(x). \end{aligned}$$

Prove that  $u$  satisfies the maximum principle

$$\max_{0 < t < T, x \in \mathbb{R}} u(t, x) \leq \max_{x \in \mathbb{R}} g(x).$$

## 4.2 Black-Scholes Equation

**Example 4.3.** Let  $f(t, S(t))$  be the price of a European put option where  $S(t)$  is the price of a stock satisfying the stochastic differential equation  $dS = \mu S dt + \sigma S dW$ , where the volatility  $\sigma$  and the drift  $\mu$  are constants. Assume also the existence of a risk free paper,  $B$ , which follows  $dB = rB dt$ , where  $r$ , the risk free rent is a constant. Find the partial differential equation of the price,  $f(t, S(t))$ , of an option.

**Solution.** Consider the portfolio  $I = -f + \alpha S + \beta B$  for  $\alpha(t), \beta(t) \in \mathbb{R}$ . Then the Itô formula and self financing, i.e.  $dI = -df + \alpha dS + \beta dB$ , imply

$$\begin{aligned} dI &= -df + \alpha dS + \beta dB \\ &= -(f_t + \mu S f_s + \frac{1}{2}\sigma^2 S^2 f_{ss})dt - f_s \sigma S dW + \alpha(\mu S dt + \sigma S dW) + \beta r B dt \\ &= \left( -(f_t + \mu S f_s + \frac{1}{2}\sigma^2 S^2 f_{ss}) + (\alpha \mu S + \beta r B) \right) dt + (-f_s + \alpha) \sigma S dW. \end{aligned}$$

Now choose  $\alpha$  such that the portfolio  $I$  becomes riskless, i.e.  $\alpha = f_s$ , so that

$$\begin{aligned} dI &= \left( -(f_t + \mu S f_s + \frac{1}{2}\sigma^2 S^2 f_{ss}) + (f_s \mu S + \beta r B) \right) dt \\ &= \left( -(f_t + \frac{1}{2}\sigma^2 S^2 f_{ss}) + \beta r B \right) dt. \end{aligned} \tag{4.2}$$

Assume also that the existence of an arbitrage opportunity is precluded, i.e.  $dI = rI dt$ , where  $r$  is the interest rate for riskless investments, to obtain

$$\begin{aligned} dI &= r(-f + \alpha S + \beta B) dt \\ &= r(-f + f_s S + \beta B) dt. \end{aligned} \tag{4.3}$$



Equation (4.2) and (4.3) show that

$$f_t + rsf_s + \frac{1}{2}\sigma^2 s^2 f_{ss} = rf, \quad t < T, \quad (4.4)$$

and finally at the maturity time  $T$  the contract value is given by definition, e.g. a standard European put option satisfies for a given exercise price  $K$

$$f(T, s) = \max(K - s, 0).$$

The deterministic partial differential equation (4.4) is called the Black-Scholes equation. The existence of adapted  $\beta$  is shown in the exercise below.  $\square$

**Exercise 4.4** (Replicating portfolio). It is said that the self financing portfolio,  $\alpha S + \beta B$ , replicates the option  $f$ . Show that there exists an adapted stochastic process  $\beta(t)$ , satisfying self financing,  $d(\alpha S + \beta B) = \alpha dS + \beta dB$ , with  $\alpha = f_s$ .

**Exercise 4.5.** Verify that the corresponding equation (4.4) holds if  $\mu, \sigma$  and  $r$  are given functions of time and stock price.

**Exercise 4.6** (Simulation of a replicating portfolio). Assume that the previously described Black-Scholes model holds and consider the case of a bank that has written (sold) a call option on the stock  $S$  with the parameters

$$S(0) = S_0 = 760, \quad r = 0.06, \quad \sigma = 0.65, \quad K = S_0.$$

with an exercise date,  $T = 1/4$  years. The goal of this exercise is to simulate the replication procedure described in Exercise 4.4, using the exact solution of the Black-Scholes call price, computed by the Octave/Matlab code

```
% Black-Scholes call option computation
function y = bsch(S,T,K,r,sigma);

normal = inline('(1+erf(x/sqrt(2)))/2','x');
d1 = (log(S/K)+(r+.5*sigma^2)*T)/sigma/sqrt(T);
d2 = (log(S/K)+(r-.5*sigma^2)*T)/sigma/sqrt(T);
y = S*normal(d1)-K*exp(-r*T)*normal(d2);
```

To this end, choose a number of hedging dates,  $N$ , and time steps  $\Delta t \equiv T/N$ . Assume that  $\beta(0) = -f_S(0, S_0)$  and then

- Write a code that computes the  $\Delta \equiv \partial f(0, S_0)/\partial S_0$  of a call option.
- Generate a realization for  $S(n\Delta t, \omega)$ ,  $n = 0, \dots, N$ .
- Generate the corresponding time discrete realizations for the processes  $\alpha_n$  and  $\beta_n$  and the portfolio value,  $\alpha_n S_n + \beta_n B_n$ .
- Generate the value after settling the contract at time  $T$ ,

$$\alpha_N S_N + \beta_N B_N - \max(S_N - K, 0).$$

Compute with only one realization, and several values of  $N$ , say  $N = 10, 20, 40, 80$ . What do you observe? How would you proceed if you don't have the exact solution of the Black-Scholes equation?

**Theorem 4.7** (Feynman-K ac). *Suppose that  $a, b, g, h$  and  $V$  are bounded smooth functions. Let  $X$  be the solution of the stochastic differential equation  $dX(t) = a(t, X(t))dt + b(t, X(t))dW(t)$  and let*

$$\begin{aligned} u(x, t) &= E[g(X(T))e^{\int_t^T V(s, X(s))ds} | X(t) = x] \\ &+ E\left[-\int_t^T h(s, X(s))e^{\int_t^s V(\tau, X(\tau))d\tau} ds | X(t) = x\right]. \end{aligned}$$

Then  $u$  is the solution of the partial differential equation

$$\begin{aligned} L_V^* u &\equiv u_t + au_x + \frac{1}{2}b^2 u_{xx} + Vu = h, \quad t < T \\ u(x, T) &= g(x). \end{aligned} \quad (4.5)$$

*Proof.* Define  $\hat{u}$  to be the solution of the equation (4.5), i.e.  $L_V^* \hat{u} = h$  and let  $G(s) \equiv e^{\int_t^s V(\tau, X(\tau)) d\tau}$ . We want to verify that  $\hat{u}$  is the claimed expected value. We have by Itô's formula, with  $L^* \hat{u} = \hat{u}_t + a\hat{u}_x + \frac{1}{2}b^2 \hat{u}_{xx}$ ,

$$\begin{aligned} d(\hat{u}(s, X(s))e^{\int_t^s V(\tau, X(\tau)) d\tau}) &= d(\hat{u}(s, X(s))G) \\ &= Gd\hat{u} + \hat{u}dG \\ &= G(L^* \hat{u} dt + b\hat{u}_x dW) + \hat{u}VG dt, \end{aligned}$$

Integrate both sides from  $t$  to  $T$ , take the expected value and use  $L^* \hat{u} = L_V^* \hat{u} - V\hat{u} = h - V\hat{u}$  to obtain

$$\begin{aligned} E[g(X(T))G(T) \mid X(t) = x] - \hat{u}(x, t) &= E\left[\int_t^T GL^* \hat{u} ds\right] + E\left[\int_t^T bG\hat{u}_x dW\right] + E\left[\int_t^T \hat{u}VG ds\right] \\ &= E\left[\int_t^T hG ds\right] - E\left[\int_t^T \hat{u}VG ds\right] + E\left[\int_t^T \hat{u}VG ds\right] \\ &= E\left[\int_t^T hG ds \mid X(t) = x\right]. \end{aligned}$$

Therefore

$$\hat{u}(x, t) = E[g(X(T))G(T) \mid X(t) = x] - E\left[\int_t^T hG ds \mid X(t) = x\right].$$

□

**Remark 4.8.** Compare Black-Scholes equation (4.4) with Equation (4.5): then  $u$  corresponds to  $f$ ,  $X$  to  $\tilde{S}$ ,  $a(t, x) = rx$ ,  $b(t, x) = \sigma x$ ,  $V = -r$  and  $h = 0$ . Using the Feynman-Kac formula, we obtain

$f(t, \tilde{S}(t)) = E[e^{-r(T-t)} \max(K - \tilde{S}(T), 0)]$ , with  $d\tilde{S} = r\tilde{S}dt + \sigma\tilde{S}dW$ , which establishes the important relation between approximation based on the Monte Carlo method and partial differential equations discussed in Chapter 1.

**Corollary 4.9.** Let  $u(x, t) = E[g(X(T)) \mid X(t) = x] = \int_{\mathbb{R}} g(y)P(y, T; x, t) dy$ . Then the density,  $P$  as a function of the first two variables, solves the Kolmogorov forward equation, also called the Fokker-Planck equation,

$$\underbrace{-\partial_s P(y, s; x, t) - \partial_y(a(y, s)P(y, s; x, t)) + \frac{1}{2}\partial_y^2(b^2(y, s)P(y, s; x, t))}_{=:LP} = 0, \quad s > t$$

$$P(y, t; x, t) = \delta(x - y),$$

where  $\delta$  is the Dirac-delta measure concentrated at zero.

*Proof.* Assume  $L\hat{P} = 0$ ,  $\hat{P}(y, t; x, t) = \delta(x - y)$ . The Feynman-Kac formula implies  $L^*u = 0$ , so that integration by part shows

$$\begin{aligned} 0 &= \int_t^T \int_{\mathbb{R}} L_{y,s}^* u(y, s) \hat{P}(y, s; x, t) dy ds \\ &= \left[ \int_{\mathbb{R}} u(y, s) \hat{P}(y, s; x, t) dy \right]_{s=t}^{s=T} + \int_t^T \int_{\mathbb{R}} u(y, s) L_{y,s} \hat{P}(y, s; x, t) dy ds \\ &= \left[ \int_{\mathbb{R}} u(y, s) \hat{P}(y, s; x, t) dy \right]_{s=t}^{s=T}. \end{aligned}$$

Consequently,

$$\begin{aligned} u(x, t) &= \int_{\mathbb{R}} g(y) \hat{P}(y, T; x, t) dy \\ &= E[g(X(T)) | X(t) = x], \end{aligned}$$

for all functions  $g$ . Therefore  $\hat{P}$  is the density function  $P$ . Hence  $P$  solves  $LP = 0$ .  $\square$

**Exercise 4.10** (Limit probability distribution). Consider the Ornstein-Uhlenbeck process defined by

$$\begin{aligned} dX(s) &= (m - X(s))ds + \sqrt{2}dW(s), \\ X(0) &= x_0. \end{aligned}$$

Verify by means of the Fokker-Plank equation that there exist a limit distribution for  $X(s)$ , when  $s \rightarrow \infty$ .

**Exercise 4.11.** Assume that  $S(t)$  is the price of a single stock. Derive a Monte-Carlo and a PDE method to determine the price of a contingent claim with the contract  $\int_0^T h(t, S(t)) dt$ , for a given function  $h$ , replacing the usual contract  $\max(S(T) - K, 0)$  for European call options.

**Exercise 4.12.** Derive the Black-Scholes equation for a general system of stocks  $S(t) \in \mathbb{R}^d$  solving

$$dS_i = a_i(t, S(t))dt + \sum_{j=1}^d b_{ij}(t, S(t))dW_j(t)$$

and a rainbow option with the contract  $f(T, S(T)) = g(S(T))$  for a given function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ , for example

$$g(S) = \max \left( \frac{1}{d} \sum_{i=1}^d S_i - K, 0 \right).$$

# Chapter 5

## The Monte-Carlo Method

This chapter gives the basic understanding of simulation of expected values  $E[g(X(T))]$  for a solution,  $X$ , of a given stochastic differential equation with a given function  $g$ . In general the approximation error has the two parts of statistical error and time discretization error, which are analyzed in the next sections. The estimation of statistical error is based on the Central Limit Theorem. The error estimate for the time discretization error of the Euler method is directly related to the proof of Feynman-Kac's theorem with an additional residual term measuring the accuracy of the approximation, which turns out to be first order in contrast to the half order accuracy for strong approximation.

### 5.1 Statistical Error

Consider the stochastic differential equation

$$dX(t) = a(t, X(t))dt + b(t, X(t))dW(t)$$

on  $t_0 \leq t \leq T$ , how can one compute the value  $E[g(X(T))]$ ? The Monte-Carlo method is based on the approximation

$$E[g(X(T))] \simeq \sum_{j=1}^N \frac{g(\bar{X}(T; \omega_j))}{N},$$

where  $\bar{X}$  is an approximation of  $X$ , e.g. the Euler method. The error in the Monte-Carlo method is

$$\begin{aligned} E[g(X(T))] - \sum_{j=1}^N \frac{g(\bar{X}(T; \omega_j))}{N} \\ = E[g(X(T)) - g(\bar{X}(T))] - \sum_{j=1}^N \frac{g(\bar{X}(T; \omega_j)) - E[g(\bar{X}(T))]}{N}. \end{aligned} \quad (5.1)$$

In the right hand side of the error representation (5.1), the first part is the time discretization error, which we will consider in the next subsection, and the second part is the statistical error, which we study here.

**Example 5.1.** Compute the integral  $I = \int_{[0,1]^d} f(x)dx$  by the Monte Carlo method, where we assume  $f(x) : [0, 1]^d \rightarrow \mathbf{R}$ .

**Solution.** We have

$$\begin{aligned}
 I &= \int_{[0,1]^d} f(x) dx \\
 &= \int_{[0,1]^d} f(x)p(x) dx \quad (\text{where } p \text{ is the uniform density function}) \\
 &= E[f(x)] \quad (\text{where } x \text{ is uniformly distributed in } [0, 1]^d) \\
 &\simeq \sum_{n=1}^N \frac{f(x(\omega_n))}{N} \\
 &\equiv I_N,
 \end{aligned}$$

where  $\{x(\omega_n)\}$  is sampled uniformly in the cube  $[0, 1]^d$ , by sampling the components  $x_i(\omega_n)$  independent and uniformly on the interval  $[0, 1]$ .  $\square$

The Central Limit Theorem is the fundamental result to understand the statistical error of Monte Carlo methods.

**Theorem 5.2** (The Central Limit Theorem). *Assume  $\xi_n$ ,  $n = 1, 2, 3, \dots$  are independent, identically distributed (i.i.d) and  $E[\xi_n] = 0$ ,  $E[\xi_n^2] = 1$ . Then*

$$\sum_{n=1}^N \frac{\xi_n}{\sqrt{N}} \rightarrow \nu, \tag{5.2}$$

where  $\nu$  is  $N(0, 1)$  and  $\rightarrow$  denotes convergence of the distributions, also called weak convergence, i.e. the convergence (5.2) means  $E[g(\sum_{n=1}^N \xi_n/\sqrt{N})] \rightarrow E[g(\nu)]$  for all bounded and continuous functions  $g$ .

*Proof.* Let  $f(t) = E[e^{it\xi_n}]$ . Then

$$f^{(m)}(t) = E[i^m \xi_n^m e^{it\xi_n}], \tag{5.3}$$

and

$$\begin{aligned}
 E[e^{it \sum_{n=1}^N \xi_n/\sqrt{N}}] &= f\left(\frac{t}{\sqrt{N}}\right)^N \\
 &= \left(f(0) + \frac{t}{\sqrt{N}}f'(0) + \frac{1}{2} \frac{t^2}{N}f''(0) + o\left(\frac{t^2}{N}\right)\right)^N.
 \end{aligned}$$

The representation (5.3) implies

$$\begin{aligned}
 f(0) &= E[1] = 1, \\
 f'(0) &= iE[\xi_n] = 0, \\
 f''(0) &= -E[\xi_n^2] = -1.
 \end{aligned}$$

Therefore

$$\begin{aligned}
E[e^{it\sum_{n=1}^N \xi_n/\sqrt{N}}] &= \left(1 - \frac{t^2}{2N} + o\left(\frac{t^2}{N}\right)\right)^N \\
&\rightarrow e^{-t^2/2}, \quad \text{as } N \rightarrow \infty \\
&= \int_{\mathbb{R}} \frac{e^{itx} e^{-x^2/2}}{\sqrt{2\pi}} dx,
\end{aligned} \tag{5.4}$$

and we conclude that the Fourier transform (i.e. the characteristic function) of  $\sum_{n=1}^N \xi_n/\sqrt{N}$  converges to the right limit of Fourier transform of the standard normal distribution. It is a fact, cf. [D], that convergence of the Fourier transform together with continuity of the limit Fourier transform at 0 implies weak convergence, so that  $\sum_{n=1}^N \xi_n/\sqrt{N} \rightharpoonup \nu$ , where  $\nu$  is  $N(0, 1)$ . The exercise below verifies this last conclusion, without reference to other results.  $\square$

**Exercise 5.3.** Show that (5.4) implies

$$E\left[g\left(\sum_{n=1}^N \xi_n/\sqrt{N}\right)\right] \rightarrow E[g(\nu)] \tag{5.5}$$

for all bounded continuous functions  $g$ . Hint: study first smooth and quickly decaying functions  $g_s$ , satisfying  $g_s(x) = \int_{-\infty}^{\infty} e^{-itx} \hat{g}_s(t) dt / (2\pi)$  with the Fourier transform  $\hat{g}_s$  of  $g_s$  satisfying  $\hat{g}_s \in L^1(\mathbb{R})$ ; show that (5.4) implies

$$E\left[g_s\left(\sum_{n=1}^N \xi_n/\sqrt{N}\right)\right] \rightarrow E[g_s(\nu)];$$

then use Chebychevs inequality to verify that no mass of  $\sum_{n=1}^N \xi_n/\sqrt{N}$  escapes to infinity; finally, let  $\chi(x)$  be a smooth cut-off function which is one for  $|x| \leq N$  and zero for  $|x| > 2N$  and split the general bounded continuous function  $g$  into  $g = g_s + g(1 - \chi) + (g\chi - g_s)$ , where  $g_s$  is an arbitrary close approximation to  $g\chi$ ; use the conclusions above to prove (5.5).

**Example 5.4.** What is the error of  $I_N - I$  in Example 5.1?

**Solution.** Let the error  $\epsilon_N$  be defined by

$$\begin{aligned}
\epsilon_N &= \sum_{n=1}^N \frac{f(x_n)}{N} - \int_{[0,1]^d} f(x) dx \\
&= \sum_{n=1}^N \frac{f(x_n) - E[f(x)]}{N}.
\end{aligned}$$

By the Central Limit Theorem,  $\sqrt{N}\epsilon_N \rightarrow \sigma\nu$ , where  $\nu$  is  $N(0,1)$  and

$$\begin{aligned}\sigma^2 &= \int_{[0,1]^d} f^2(x)dx - \left( \int_{[0,1]^d} f(x)dx \right)^2 \\ &= \int_{[0,1]^d} \left( f(x) - \int_{[0,1]^d} f(x)dx \right)^2 dx.\end{aligned}$$

In practice,  $\sigma^2$  is approximated by

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{n=1}^N \left( f(x_n) - \sum_{m=1}^N \frac{f(x_m)}{N} \right)^2.$$

□

One can generate approximate random numbers, so called pseudo random numbers, by for example the method

$$\xi_{i+1} \equiv a\xi_i + b \pmod{n}$$

where  $a$  and  $n$  are relative prime and the initial  $\xi_0$  is called the seed, which determines all other  $\xi_i$ . For example the combinations  $n = 2^{31}$ ,  $a = 2^{16} + 3$  and  $b = 0$ , or  $n = 2^{31} - 1$ ,  $a = 7^5$  and  $b = 0$  are used in practise. In Monte Carlo computations, we use the pseudo random numbers  $\{x_i\}_{i=1}^N$ , where  $x_i = \frac{\xi_i}{n} \in [0, 1]$ , which for  $N \ll 2^{31}$  behave approximately as independent uniformly distributed variables.

**Theorem 5.5.** *The following Box-Müller method generates two independent normal random variables  $x_1$  and  $x_2$  from two independent uniformly distributed variables  $y_1$  and  $y_2$*

$$\begin{aligned}x_1 &= \sqrt{-2 \log(y_2)} \cos(2\pi y_1) \\ x_2 &= \sqrt{-2 \log(y_2)} \sin(2\pi y_1).\end{aligned}$$

**Sketch of the Idea.** The variables  $x$  and  $y$  are independent standard normal variables if and only if their joint density function is  $e^{-(x^2+y^2)/2}/2\pi$ . We have

$$e^{-(x^2+y^2)/2} dx dy = r e^{-r^2/2} dr d\theta = d(e^{-r^2/2}) d\theta$$

using  $x = r \cos \theta$ ,  $y = r \sin \theta$  and  $0 \leq \theta < 2\pi$ ,  $0 \leq r < \infty$ . The random variables  $\theta$  and  $r$  can be sampled by taking  $\theta$  to be uniformly distributed in the interval  $[0, 2\pi)$  and  $e^{-r^2/2}$  to be uniformly distributed in  $(0, 1]$ , i.e.  $\theta = 2\pi y_1$ , and  $r = \sqrt{-2 \log(y_2)}$ . □

**Example 5.6.** Consider the stochastic differential equation  $dS = rSdt + \sigma SdW$ , in the risk neutral formulation where  $r$  is the riskless rate of return and  $\sigma$  is the volatility. Then

$$S_T = S_0 e^{rT - \frac{\sigma^2}{2}T + \sigma\sqrt{T}\nu}$$

where  $\nu$  is  $N(0,1)$ . The values of a call option,  $f_c$ , and put option,  $f_p$ , are by Remark 4.8

$$f_c = e^{-rT} E[\max(S(T) - K, 0)]$$

and

$$f_p = e^{-rT} E[\max(K - S(T), 0)].$$



**Example 5.7.** Consider the system of stochastic differential equations,

$$dS_i = rS_i dt + \sum_{j=1}^M \sigma_{ij} S_i dW_j, \quad i = 1, \dots, M.$$

Then

$$S_i(T) = S_i(0) e^{rT - \sum_{j=1}^M \left( \sigma_{ij} \sqrt{T} \nu_j - \frac{\sigma_{ij}^2}{2} T \right)}$$

where  $\nu_j$  are independent and  $N(0, 1)$ . A rainbow call option, based on  $S_{av} = \frac{1}{M} \sum_{i=1}^M S_i$ , can then be simulated by the Monte Carlo method and

$$f_c = e^{-rT} E[\max(S_{av}(T) - K, 0)].$$

## 5.2 Time Discretization Error

Consider the stochastic differential equation

$$dX(t) = a(t, X(t))dt + b(t, X(t))dW(t), \quad 0 \leq t \leq T,$$

and let  $\bar{X}$  be the forward Euler discretization of  $X$ . Then

$$\bar{X}(t_{n+1}) - \bar{X}(t_n) = a(t_n, \bar{X}(t_n))\Delta t_n + b(t_n, \bar{X}(t_n))\Delta W_n, \quad (5.6)$$

where  $\Delta t_n = t_{n+1} - t_n$  and  $\Delta W_n = W(t_{n+1}) - W(t_n)$  for a given discretization  $0 = t_0 < t_1 < \dots < t_N = T$ . Equation (5.6) can be extended, for theoretical use, to all  $t$  by

$$\bar{X}(t) - \bar{X}(t_n) = \int_{t_n}^t \bar{a}(s, \bar{X}) ds + \int_{t_n}^t \bar{b}(s, \bar{X}) dW(s), \quad t_n \leq t < t_{n+1},$$

where, for  $t_n \leq s < t_{n+1}$ ,

$$\begin{aligned} \bar{a}(s, \bar{X}) &= a(t_n, \bar{X}(t_n)), \\ \bar{b}(s, \bar{X}) &= b(t_n, \bar{X}(t_n)). \end{aligned} \quad (5.7)$$

**Theorem 5.8.** Assume that  $a, b$  and  $g$  are smooth and decay sufficiently fast as  $|x| \rightarrow \infty$ . Then there holds

$$E[g(X(T)) - g(\bar{X}(T))] = \mathcal{O}(\max \Delta t).$$

*Proof.* Let  $u$  satisfy the equation

$$L^* u \equiv u_t + au_x + \frac{b^2}{2} u_{xx} = 0, \quad t < T \quad (5.8)$$

$$u(x, T) = g(x). \quad (5.9)$$

The Feynman-Kác formula shows

$$u(x, t) = E[g(X(T)) | X(t) = x]$$

and in particular

$$u(0, X(0)) = E[g(X(T))]. \quad (5.10)$$

Then by the Itô formula,

$$\begin{aligned} du(t, \bar{X}(t)) &= \left( u_t + \bar{a}u_x + \frac{\bar{b}^2}{2}u_{xx} \right) (t, \bar{X}(t))dt + \bar{b}u_x(t, \bar{X}(t))dW \\ &\stackrel{(5.8)}{=} \left( -au_x - \frac{b^2}{2}u_{xx} + \bar{a}u_x + \frac{\bar{b}^2}{2}u_{xx} \right) (t, \bar{X}(t))dt + \bar{b}u_x(t, \bar{X}(t))dW \\ &= \left\{ (\bar{a} - a)u_x(t, \bar{X}(t)) + \left( \frac{\bar{b}^2}{2} - \frac{b^2}{2} \right) u_{xx}(t, \bar{X}(t)) \right\} dt \\ &\quad + \bar{b}(t, \bar{X})u_x(t, \bar{X}(t))dW. \end{aligned}$$

Evaluate the integral from 0 to T,

$$\begin{aligned} u(T, \bar{X}(T)) - u(0, X(0)) &= \int_0^T (\bar{a} - a)u_x(t, \bar{X}(t))dt + \int_0^T \frac{\bar{b}^2 - b^2}{2}u_{xx}(t, \bar{X}(t))dt \\ &\quad + \int_0^T \bar{b}(t, \bar{X}(t))u_x dW. \end{aligned}$$

Take the expected value and use (5.10) to obtain

$$\begin{aligned} E[g(\bar{X}(T)) - g(X(T))] &= \int_0^T E[(\bar{a} - a)u_x] + \frac{1}{2}E[(\bar{b}^2 - b^2)u_{xx}]dt + E \left[ \int_0^T \bar{b}u_x dW \right] \\ &= \int_0^T E[(\bar{a} - a)u_x] + \frac{1}{2}E[(\bar{b}^2 - b^2)u_{xx}]dt. \end{aligned}$$

The following Lemma 5.9 proves the Theorem. □

**Lemma 5.9.** *There holds for  $t_n \leq t < t_{n+1}$*

$$\begin{aligned} f_1(t) &\equiv E[(\bar{a}(t, \bar{X}) - a(t, \bar{X}(t)))u_x(t, \bar{X}(t))] = \mathcal{O}(\Delta t_n), \\ f_2(t) &\equiv E[(\bar{b}^2(t, \bar{X}) - b^2(t, \bar{X}(t)))u_{xx}(t, \bar{X}(t))] = \mathcal{O}(\Delta t_n). \end{aligned}$$

*Proof.* Since  $\bar{a}(t, \bar{X}) = a(t_n, \bar{X}(t_n))$ ,

$$f_1(t_n) = E[(\bar{a}(t_n, \bar{X}) - a(t_n, \bar{X}(t_n)))u_x(t_n, \bar{X}(t_n))] = 0. \quad (5.11)$$

Provided  $|f'_1(t)| \leq C$ , the initial condition (5.11) implies that  $f_1(t) = \mathcal{O}(\Delta t_n)$ , for  $t_n \leq t < t_{n+1}$ . Therefore, it remains to show that  $|f'_1(t)| \leq C$ . Let  $\alpha(t, x) = -(a(t, x) -$

$a(t_n, \bar{X}(t_n))u_x(t, x)$ , so that  $f(t) = E[\alpha(t, \bar{X}(t))]$ . Then by Itô's formula

$$\begin{aligned} \frac{df}{dt} &= \frac{d}{dt} E[\alpha(t, \bar{X}(t))] = E[d\alpha(t, \bar{X}(t))] / dt \\ &= E\left[\left(\alpha_t + \bar{a}\alpha_x + \frac{\bar{b}^2}{2}\alpha_{xx}\right) dt + \alpha_x \bar{b}dW\right] / dt \\ &= E\left[\alpha_t + \bar{a}\alpha_x + \frac{\bar{b}^2}{2}\alpha_{xx}\right] \\ &= \mathcal{O}(1). \end{aligned}$$

Therefore there exists a constant  $C$  such that  $|f'(t)| \leq C$ , for  $t_n < t < t_{n+1}$ , and consequently

$$f_1(t) \equiv E[(\bar{a}(t, \bar{X}) - a(t, \bar{X}(t)))u_x(t, \bar{X}_t)] = \mathcal{O}(\Delta t_n), \quad \text{for } t_n \leq t < t_{n+1}.$$

Similarly, we can also prove

$$f_2(t) \equiv E[(\bar{b}^2(t, \bar{X}) - b^2(t, \bar{X}(t)))u_{xx}(t, \bar{X}_t)] = \mathcal{O}(\Delta t_n), \quad \text{for } t_n \leq t < t_{n+1}.$$

□

**Example 5.10.** Consider the stochastic volatility model,

$$\begin{aligned} dS &= \omega S dt + \sigma S dZ \\ d\sigma &= \alpha \sigma dt + \nu \sigma dW \end{aligned} \tag{5.12}$$

where  $Z$  and  $W$  are Brownian motions with correlation coefficient  $\rho$ , i.e.  $E[dZdW] = \rho dt$ . We can then construct  $Z$  and  $W$  from the independent  $W_1$  and  $W_2$  by

$$W = W_1, \quad Z = \rho W_1 + \sqrt{1 - \rho^2} W_2.$$

**Exercise 5.11.** In the risk neutral formulation a stock price solves the stochastic differential equation

$$dS = rSdt + \sigma SdW(t),$$

with constant interest rate  $r$  and volatility  $\sigma$ .

(i) Show that

$$S(T) = S(0)e^{rT - \frac{\sigma^2}{2}T + \sigma W(T)}. \quad (5.13)$$

(ii) Use equation (5.13) to simulate the price

$$f(0, S(0)) = e^{-rT} E[ \max (S(T) - K, 0) ]$$

of an European call option by a Monte-Carlo method.

(iii) Compute also the corresponding  $\Delta = \partial f(0, S)/\partial S$  by approximating with a difference quotient and determine a good choice of your approximation of " $\partial S$ ".

(iv) Estimate the accuracy of your results. Suggest a better method to solve this problem.

**Exercise 5.12.** Assume that a system of stocks solves

$$\frac{dS_i}{S_i(t)} = rdt + \sum_{j=1}^d \sigma_{ij} dW_j(t) \quad i = 1, \dots, d$$

where  $W_j$  are independent Brownian motions.

(i) Show that

$$S_i(T) = S(0)e^{rT + \sum_{j=1}^d (\sigma_{ij} W_j(T) - \frac{1}{2} \sigma_{ij}^2 T)}.$$

(ii) Let  $S_{av} \equiv \sum_{i=1}^d S_i/d$  and simulate the price of the option above with  $S(T)$  replaced by  $S_{av}(T)$ . Estimate the accuracy of your results. Can you find a better method to solve this problem?

**Exercise 5.13** (An example of variance reduction). Consider the computation of a call option on an index  $Z$ ,

$$\pi_t = e^{-r(T-t)} E[\max(Z(T) - K, 0)], \quad (5.14)$$

where  $Z$  is the average of  $d$  stocks,

$$Z(t) \equiv \frac{1}{d} \sum_{i=1}^d S_i(t)$$

and

$$dS_i(t) = rS_i(t)dt + \sigma_i S_i(t)dW_i(t), \quad i = 1, \dots, d$$

with volatilities

$$\sigma_i \equiv 0.2 * (2 + \sin(i)) \quad i = 1, \dots, d.$$

The correlation between Wiener processes is given by

$$E[dW_i(t)dW_{i'}(t)] = \exp(-2|i - i'|/d)dt \quad 1 \leq i, i' \leq d.$$

The goal of this exercise is to experiment with two different variance reduction techniques, namely the antithetic variates and the control variates.

From now on we take  $d = 10$ ,  $r = 0.04$  and  $T = 0.5$  in the example above.

- (i) Implement a Monte Carlo approximation with for the value in (5.14). Estimate the statistical error. Choose a number of realizations such that the estimate for the statistical error is less than 1% of the value we want to approximate.
- (ii) Same as (i) but using antithetic variates. The so called *antithetic variates* technique reduces the variance in a sample estimator  $\mathcal{A}(M; Y)$  by using another estimator  $\mathcal{A}(M; Y')$  with the same expectation as the first one, but which is negatively correlated with the first. Then, the improved estimator is  $\mathcal{A}(M; \frac{1}{2}(Y + Y'))$ . Here, the choice of  $Y$  and  $Y'$  relates to the Wiener process  $W$  and its reflection along the time axis,  $-W$ , which is also a Wiener process, i.e.

$$\pi_t \approx \frac{1}{M} \sum_{j=1}^M \frac{\{\max(Z(W(T, \omega_j)) - K, 0) + \max(Z(-W(T, \omega_j)) - K, 0)\}}{2}.$$

- (iii) Same as (i) but using control variates to reduce the variance. The control variates technique is based on the knowledge of an estimator  $Y''$ , positively correlated with  $Y$ , whose expected value  $E[Y'']$  is known and relatively close to the desired  $E[Y]$ , yielding  $Y - Y'' + E[Y'']$  as an improved estimator.

For the application of control variates to (5.14) use the geometric average

$$\hat{Z}(t) \equiv \left\{ \prod_{i=1}^d S_i(t) \right\}^{\frac{1}{d}},$$

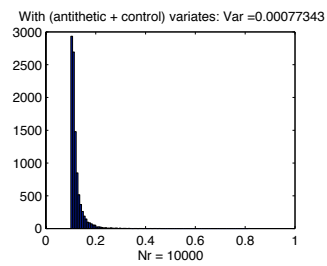
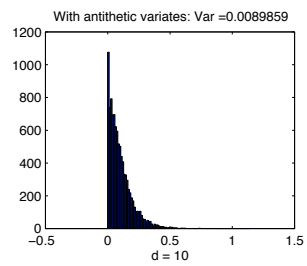
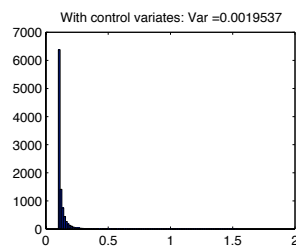
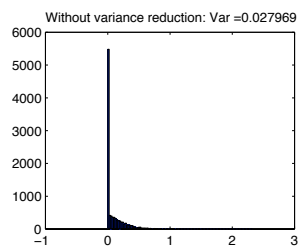
compute

$$\hat{\pi}_t = e^{-r(T-t)} E[\max(\hat{Z}(T) - K, 0)]$$

exactly (hint: find a way to apply Black-Scholes formula). Then approximate

$$\pi_t \approx \hat{\pi}_t + \frac{e^{-r(T-t)}}{M} \sum_{j=1}^M \left\{ \max(Z(W(T, \omega_j)) - K, 0) - \max(\hat{Z}(W(T, \omega_j)) - K, 0) \right\}.$$

- (iv) Discuss the results from (i)-(iii). Does it pay to use variance reduction?



# Chapter 6

## Finite Difference Methods

This section introduces finite difference methods for approximation of partial differential equations. We first apply the finite difference method to a partial differential equation for a financial option problem, which is more efficiently computed by partial differential methods than Monte Carlo techniques. Then we discuss the fundamental Lax Equivalence Theorem, which gives the basic understanding of accuracy and stability for approximation of differential equations.

### 6.1 American Options

Assume that the stock value,  $S(t)$ , evolves in the risk neutral formulation by the Itô geometric Brownian motion

$$dS = rSdt + \sigma SdW.$$

An American put option is a contract that gives the possibility to sell a stock for a fixed price  $K$  up to time  $T$ . Therefore the derivation of option values in Chapter 4 shows that European and American options have the formulations:

- (i) The price of an European put option is

$$f(t, s) \equiv E[ e^{-r(T-t)} \max(K - S(T), 0) | S(t) = s ].$$

- (ii) The price of an American option is obtained by maximizing over all sell time  $\tau$  strategies, which depend on the stock price up to the sell time,

$$f_A(t, s) \equiv \max_{t \leq \tau \leq T} E[ e^{-r(\tau-t)} \max(K - S(\tau), 0) | S(t) = s ]. \quad (6.1)$$

How to find the optimal selling strategy for an American option? Assume that selling is only allowed at the discrete time levels  $0, \Delta t, 2\Delta t, \dots, T$ . Consider the small time step  $(T - \Delta t, T)$ . By assumption the option is not sold in the step. Therefore the European value  $f(t, s)$  holds, where  $f(T, s) = \max(K - s, 0)$  and for  $T - \Delta t < t < T$

$$f_t + rSf_S + \frac{1}{2}\sigma^2 S^2 f_{SS} = rf. \quad (6.2)$$

If, for a fixed stock price  $s = S(T - \Delta t)$ , there holds  $f(T - \Delta t, s) < \max(K - s, 0)$  then keeping the option gives the expected value  $f(T - \Delta t, s)$  which is clearly less than the value  $\max(K - s, 0)$  obtained by selling at time  $T - \Delta t$ . Therefore it is optimal to sell if  $f(T - \Delta t, s) < \max(K - s, 0) \equiv f_F$ . Modify the initial data at  $t = T - \Delta t$  to  $\max(f(T - \Delta t, s), f_F)$  and repeat the step (6.2) for  $(T - 2\Delta t, T - \Delta t)$  and so on. The price of the American option is obtained as the limit of this solution as  $\Delta t \rightarrow 0$ .

**Example 6.1.** A corresponding Monte Carlo method based on (6.1) requires simulation of expected values  $E[e^{-r\tau} \max(K - S(\tau), 0)]$  for many different possible selling time strategies  $\tau$  until an approximation of the maximum values is found. Since the  $\tau$  need to depend on  $\omega$ , with  $M$  time steps and  $N$  realizations there are  $M^N$  different strategies.

Note that the optimal selling strategy

$$\tau = \tau^* = \inf_v \{v : t \leq v \leq T, f_A(v, S(v)) = \max(K - S(v), 0)\}$$

for the American option, which is a function of  $f_A$ , seems expensive to evaluate by Monte Carlo technique, but is obtained directly in the partial differential formulation above and below. This technique is a special case of the so called dynamic programming method, which we shall study systematically for general optimization problems in a later Chapter, cf. also the last example in Chapter 1.

Here and in Exercise 6.2 is a numerical method to determine the value of an American option:

- (1) Discretize the computational domain  $[0, T] \times [s_0, s_1]$  and let

$$f_A(n\Delta t, i\Delta S) \simeq \bar{f}_{n,i}, \quad \bar{f}_{N,i} = \max(K - i\Delta S, 0).$$

- (2) Use the Euler and central difference methods for the equation (6.2)

$$\begin{aligned} \partial_t f_A &\simeq \frac{\bar{f}_{n,i} - \hat{f}_{n-1,i}}{\Delta t} & \partial_S f_A &\simeq \frac{\bar{f}_{n,i+1} - \bar{f}_{n,i-1}}{2\Delta S} \\ \partial_{SS} f_A &\simeq \frac{\bar{f}_{n,i+1} - 2\bar{f}_{n,i} + \bar{f}_{n,i-1}}{(\Delta S)^2} & f_A &\simeq \bar{f}_{n,i}. \end{aligned}$$

- (3) Make a Black-Scholes prediction for each time step

$$\begin{aligned} \hat{f}_{n-1,i} &= \bar{f}_{n,i}(1 - r\Delta t - \sigma^2 i^2 \Delta t) + \bar{f}_{n,i+1} \left( \frac{1}{2} r i \Delta t + \frac{1}{2} \sigma^2 i^2 \Delta t \right) \\ &+ \bar{f}_{n,i-1} \left( -\frac{1}{2} r i \Delta t + \frac{1}{2} \sigma^2 i^2 \Delta t \right). \end{aligned}$$

- (4) Compare the prediction with selling by letting

$$\bar{f}_{n-1,i} = \max(\hat{f}_{n-1,i}, \max(K - i\Delta S, 0)),$$

and go to the next time Step 3 by decreasing  $n$  by 1.



**Exercise 6.2.** The method above needs in addition boundary conditions at  $S = s_0$  and  $S = s_1$  for  $t < T$ . How can  $s_0, s_1$  and these conditions be chosen to yield a good approximation?

**Exercise 6.3.** Give a trinomial tree interpretation of the finite difference scheme

$$\begin{aligned}\bar{f}_{n+1,i} &= \bar{f}_{n,i}(1 + r\Delta t + \sigma^2 i^2 \Delta t) + \bar{f}_{n,i+1}\left(-\frac{1}{2}ri\Delta t - \frac{1}{2}\sigma^2 i^2 \Delta t\right) \\ &+ \bar{f}_{n,i-1}\left(\frac{1}{2}ri\Delta t - \frac{1}{2}\sigma^2 i^2 \Delta t\right),\end{aligned}$$

for Black-Scholes equation of an European option. Binomial and trinomial tree approximations are frequent in the finance economy literature, cf. [J. Hull].

Let us now study general finite difference methods for partial differential equations. The motivation to introduce general finite difference methods in contrast to study only the binomial and trinomial tree methods is that higher order methods, such as the Crank-Nicolson method below, are more efficient to solve e.g. (6.2).

The error for the binomial and the trinomial tree method applied to the partial differential equation (6.2) for a European option is  $\varepsilon = \mathcal{O}(\Delta t + (\Delta s)^2)$ , which is clearly the same for the related forward and backward Euler methods. The work is then  $\mathcal{A} = \mathcal{O}((\Delta t \Delta s)^{-1})$ , so that  $\mathcal{A} = \mathcal{O}(\varepsilon^{-3/2})$ . For the Crank-Nicolson method the accuracy is  $\varepsilon = \mathcal{O}((\Delta t)^2 + (\Delta s)^2)$  and the work is still  $\mathcal{A} = \mathcal{O}((\Delta t \Delta s)^{-1})$ , which implies the improved bound  $\mathcal{A} = \mathcal{O}(\varepsilon^{-1})$ . For a general implicit method with a smooth exact solution in  $[0, T] \times \mathbb{R}^d$  the accuracy is  $\varepsilon = \mathcal{O}((\Delta t)^q + (\Delta s)^p)$  with the minimal work ( using e.g. the multigrid method )  $\mathcal{A} = \mathcal{O}(\frac{q^2}{\Delta t} (\frac{p^2}{\Delta s})^d)$ , which gives  $\mathcal{A} = \mathcal{O}(\frac{q^2}{\varepsilon^{1/q}} (\frac{p^2}{\varepsilon^{1/p}})^d)$ . In the next section we derive these error estimates for some model problems.

## 6.2 Lax Equivalence Theorem

Lax equivalence theorem defines the basic concepts for approximation of linear well posed differential equations. Here, well posed means that the equation is solvable for data in a suitable function space and that the solution operator is bounded. We will first formally state the result without being mathematically precise with function spaces and norms. Then we present two examples with proofs based on norms and functions spaces.

The ingredients of Lax Equivalence Theorem 6.4 are:

- (0) an exact solution  $u$ , satisfying the *linear well posed equation*  $Lu = f$ , and an approximation  $u_h$ , obtained from  $L_h u_h = f_h$ ;
- (1) *stability*, the approximate solution operators  $\|L_h^{-1}\|$  are uniformly bounded in  $h$  and the exact solution operator  $\|L^{-1}\|$  is bounded;
- (2) *consistency*,  $f_h \rightarrow f$  and  $L_h u \rightarrow Lu$  as the mesh size  $h \rightarrow 0$ ; and
- (3) *convergence*,  $u_h \rightarrow u$  as the mesh size  $h \rightarrow 0$ .

**Theorem 6.4.** *The combination of stability and consistency is equivalent to convergence.*

**The idea of the proof.** To verify convergence, consider the identity

$$u - u_h = L_h^{-1} [ L_h u - L_h u_h ] \stackrel{\text{Step(0)}}{=} L_h^{-1} [ (L_h u - Lu) + (f - f_h) ].$$

Stability implies that  $L_h^{-1}$  is bounded and consistency implies that

$$L_h u - Lu \rightarrow 0 \text{ and } f - f_h \rightarrow 0,$$

and consequently the convergence holds

$$\begin{aligned} \lim_{h \rightarrow 0} (u - u_h) &= \lim_{h \rightarrow 0} L_h^{-1} [ (L_h u - Lu) + (f - f_h) ] \\ &= 0. \end{aligned}$$

Clearly, consistency is necessary for convergence. Example 6.7, below, indicates that also stability is necessary.  $\square$

Let us now more precisely consider the requirements and norms to verify stability and consistency for two concrete examples of ordinary and partial differential equations.

**Example 6.5.** Consider the forward Euler method for the ordinary differential equation

$$\begin{aligned} u'(t) &= Au(t) \quad 0 < t < 1, \\ u(0) &= u_0. \end{aligned} \tag{6.3}$$

Verify the conditions of stability and consistency in Lax Equivalence Theorem.

**Solution.** For a given partition,  $0 = t_0 < t_1 < \dots < t_N = 1$ , with  $\Delta t = t_{n+1} - t_n$ , let

$$\begin{aligned} u_{n+1} &\equiv (I + \Delta t A) u_n \\ &= G^n u_0 \quad \text{where } G = (I + \Delta t A). \end{aligned}$$

Then:

- (1) Stability means  $|G^n| + |H^n| \leq e^{Kn\Delta t}$  for some  $K$ , where  $|\cdot|$  denotes the matrix norm  $|F| \equiv \sup_{\{v \in \mathbb{R}^n: |v| \leq 1\}} |Fv|$  with the Euclidean norm  $|w| \equiv \sqrt{\sum_i w_i^2}$  in  $\mathbb{R}^n$ .
- (2) Consistency means  $|(G - H)v| \leq C(\Delta t)^{p+1}$ , where  $H = e^{\Delta t A}$  and  $p$  is the order of accuracy. In other words, the consistency error  $(G - H)v$  is the local approximation error after one time step with the same initial data  $v$ .

This stability and consistency imply the convergence

$$\begin{aligned} |u_n - u(n\Delta t)| &= |(G^n - H^n)u_0| \\ &= |(G^{n-1} + G^{n-2}H + \dots + GH^{n-2} + H^{n-1})(G - H)u_0| \\ &\leq |G^{n-1} + G^{n-2}H + \dots + GH^{n-2} + H^{n-1}| |(G - H)u_0| \\ &\leq C(\Delta t)^{p+1} n |u_0| e^{Kn\Delta t} \\ &\leq C'(\Delta t)^p, \end{aligned}$$

with the convergence rate  $\mathcal{O}(\Delta t^p)$ . For example,  $p = 1$  in case of the Euler method and  $p = 2$  in case of the trapezoidal method.  $\square$

**Example 6.6.** Consider the heat equation

$$\begin{aligned} u_t &= u_{xx} \quad t > 0, \\ u(0) &= u_0. \end{aligned} \tag{6.4}$$

Verify the stability and consistency conditions in Lax Equivalence Theorem.

**Solution.** Apply the Fourier transform to equation (6.4),

$$\hat{u}_t = -\omega^2 \hat{u}$$

so that

$$\hat{u}(t, \omega) = e^{-t\omega^2} \hat{u}_0(\omega).$$

Therefore  $\hat{H} = e^{-\Delta t \omega^2}$  is the exact solution operator for one time step, i.e.  $\hat{u}(t + \Delta t) = \hat{H}\hat{u}(t)$ . Consider the difference approximation of (6.4)

$$\frac{u_{n+1,i} - u_{n,i}}{\Delta t} = \frac{u_{n,i+1} - 2u_{n,i} + u_{n,i-1}}{\Delta x^2},$$

which shows

$$u_{n+1,i} = u_{n,i} \left( 1 - \frac{2\Delta t}{\Delta x^2} \right) + \frac{\Delta t}{\Delta x^2} (u_{n,i+1} + u_{n,i-1}),$$

where  $u_{n,i} \simeq u(n\Delta t, i\Delta x)$ . Apply the Fourier transform to obtain

$$\begin{aligned} \hat{u}_{n+1} &= \left[ \left( 1 - \frac{2\Delta t}{\Delta x^2} \right) + \frac{\Delta t}{\Delta x^2} (e^{j\Delta x\omega} + e^{-j\Delta x\omega}) \right] \hat{u}_n \\ &= \left[ 1 - 2\frac{\Delta t}{\Delta x^2} + 2\frac{\Delta t}{\Delta x^2} \cos(\Delta x\omega) \right] \hat{u}_n \\ &= \hat{G}\hat{u}_n \quad \left( \text{Let } \hat{G} \equiv 1 - 2\frac{\Delta t}{\Delta x^2} + 2\frac{\Delta t}{\Delta x^2} \cos(\Delta x\omega) \right) \\ &= \hat{G}^{n+1}\hat{u}_0. \end{aligned}$$

(i) We have

$$\begin{aligned} 2\pi \|u_n\|_{L^2}^2 &= \|\hat{u}_n\|_{L^2}^2 \quad (\text{by Parseval's formula}) \\ &= \|\hat{G}^n \hat{u}_0\|_{L^2}^2 \\ &\leq \sup_{\omega} |\hat{G}^n|^2 \|\hat{u}_0\|_{L^2}^2. \end{aligned}$$

Therefore the condition

$$\|\hat{G}^n\|_{L^\infty} \leq e^{Kn\Delta t} \tag{6.5}$$

implies  $L^2$ -stability.

(ii) We have

$$2\pi \|u_1 - u(\Delta t)\|_{L^2}^2 = \|\hat{G}\hat{u}_0 - \hat{H}\hat{u}_0\|_{L^2}^2,$$

where  $u_1$  is the approximate solution after one time step. Let  $\lambda \equiv \frac{\Delta t}{\Delta x^2}$ , then we obtain

$$\begin{aligned} |(\hat{G} - \hat{H})\hat{u}_0| &= \left| \left(1 - 2\lambda + 2\lambda \cos \Delta x \omega - e^{-\Delta t \omega^2}\right) \hat{u}_0 \right| \\ &= \mathcal{O}(\Delta t^2) \omega^4 |\hat{u}_0|, \end{aligned}$$

since for  $0 \leq \Delta t \omega^2 \equiv x \leq 1$

$$\begin{aligned} |1 - 2\lambda + 2\lambda \cos \sqrt{x/\lambda} - e^{-x}| &= \left(1 - 2\lambda + 2\lambda \left(1 - \frac{x}{2\lambda} + \mathcal{O}(x^2)\right) - (1 - x + \mathcal{O}(x^2))\right) \\ &\leq Cx^2 = C(\Delta t)^2 \omega^4, \end{aligned}$$

and for  $1 < \Delta t \omega^2 = x$

$$|1 - 2\lambda + 2\lambda \cos \sqrt{x/\lambda} - e^{-x}| \leq C = C \frac{(\Delta t)^2 \omega^4}{x^2} \leq C(\Delta t)^2 \omega^4.$$

Therefore the consistency condition reduces to

$$\begin{aligned} \|(\hat{G} - \hat{H})\hat{u}_0\| &\leq \|K \Delta t^2 \omega^4 \hat{u}_0\| \\ &\leq K \Delta t^2 \|\partial_{xxxx} u_0\|_{L^2}. \end{aligned} \tag{6.6}$$

(iii) The stability (6.5) holds if

$$\|\hat{G}\|_{L^\infty} \equiv \sup_{\omega} |\hat{G}(\omega)| = \max_{\omega} |1 - 2\lambda + 2\lambda \cos \Delta x \omega| \leq 1, \tag{6.7}$$

which requires

$$\lambda = \frac{\Delta t}{\Delta x^2} \leq \frac{1}{2}. \tag{6.8}$$

The  $L^2$ -stability condition (6.7) is called the von Neuman stability condition.

(iv) Convergence follows by the estimates (6.6), (6.7) and  $\|\hat{H}\|_{L^\infty} \leq 1$

$$\begin{aligned} 2\pi \|u_n - u(n\Delta t)\|_{L^2}^2 &= \|(\hat{G}^n - \hat{H}^n)\hat{u}_0\|_{L^2}^2 \\ &= \|(\hat{G}^{n-1} + \hat{G}^{n-2}\hat{H} + \dots + \hat{H}^{n-1})(\hat{G} - \hat{H})\hat{u}_0\|_{L^2}^2 \\ &\leq \|\hat{G}^{n-1} + \hat{G}^{n-2}\hat{H} + \dots + \hat{H}^{n-1}\|_{L^\infty}^2 \|(\hat{G} - \hat{H})\hat{u}_0\|_{L^2}^2 \\ &\leq (Kn(\Delta t)^2)^2 \leq (KT\Delta t)^2, \end{aligned}$$

and consequently the convergence rate is  $\mathcal{O}(\Delta t)$ .  $\square$

Let us study the relations between the operators  $G$  and  $H$  for the simple model problem

$$\begin{aligned} u' + \lambda u &= 0 \\ u(0) &= 1 \end{aligned}$$

with an approximate solution  $u_{n+1} = r(x)u_n$  (where  $x = \lambda\Delta t$ ):

(1) the exact solution satisfies

$$r(x) = e^{-\lambda\Delta t} = e^{-x},$$

(2) the forward Euler method

$$\frac{u_{n+1} - u_n}{\Delta t} + \lambda u_n = 0 \Rightarrow r(x) = 1 - x,$$

(3) the backward Euler method

$$\frac{u_{n+1} - u_n}{\Delta t} + \lambda u_{n+1} = 0 \Rightarrow r(x) = (1 + x)^{-1},$$

(4) the trapezoidal method

$$\frac{u_{n+1} - u_n}{\Delta t} + \frac{\lambda}{2}(u_n + u_{n+1}) = 0 \Rightarrow r(x) = \left(1 + \frac{x}{2}\right)^{-1} \left(1 - \frac{x}{2}\right),$$

and

(5) the Lax-Wendroff method

$$u_{n+1} = u_n - \Delta t \lambda u_n + \frac{1}{2} \Delta t^2 \lambda^2 u_n \Rightarrow r(x) = 1 - x + \frac{1}{2} x^2.$$

The consistence  $|e^{-\lambda\Delta t} - r(\lambda\Delta t)| = \mathcal{O}(\Delta t^{p+1})$  holds with  $p = 1$  in case 2 and 3, and  $p = 2$  in case 4 and 5. The following stability relations hold:

- (1)  $|r(x)| \leq 1$  for  $x \geq 0$  in case 1, 3 and 4.
- (2)  $r(x) \rightarrow 0$  as  $x \rightarrow \infty$  in case 1 and 3.
- (3)  $r(x) \rightarrow 1$  as  $x \rightarrow \infty$  in case 4.

Property (1) shows that for  $\lambda > 0$  case 3 and 4 are unconditionally stable. However Property (2) and (3) refine this statement and imply that only case 3 has the same damping behavior for large  $\lambda$  as the exact solution. Although the damping Property (2) is not necessary to prove convergence it is advantageous to have for problems with many time scales, e.g. for a system of equations (6.3) where  $A$  has eigenvalues  $\lambda_i \leq 1$ ,  $i = 1, \dots, N$  and some  $\lambda_j \ll -1$ , ( why?).

The unconditionally stable methods, e.g. case 3 and 4, are in general more efficient to solve parabolic problems, such as the Black-Scholes equation (6.2), since they require for the same accuracy fewer time steps than the explicit methods, e.g. case 2 and 5. Although the work in each time step for the unconditionally stable methods may be larger than for the explicit methods.

**Exercise 6.7.** Show by an example that  $\|u_n\|_{L^2}^2 \rightarrow \infty$  if for some  $\omega$  there holds  $|\hat{G}(\omega)| > 1$ , in Example 6.6, i.e. the von Neumann stability condition does not hold.

## Chapter 7

# The Finite Element Method and Lax-Milgram's Theorem

This section presents the finite element method, including adaptive approximation and error estimates, together with the basic theory for elliptic partial differential equations. The motivation to introduce finite element methods is the computational simplicity and efficiency for construction of stable higher order discretizations for elliptic and parabolic differential equations, such as the Black and Scholes equation, including general boundary conditions and domains. Finite element methods require somewhat more work per degree of freedom as compared to finite difference methods on a uniform mesh. On the other hand, construction of higher order finite difference approximations including general boundary conditions or general domains is troublesome.

In one space dimension such an elliptic problem can, for given functions  $a, f, r : (0, 1) \rightarrow \mathbf{R}$ , take the form of the following equation for  $u : [0, 1] \rightarrow \mathbf{R}$ ,

$$\begin{aligned} (-au')' + ru &= f && \text{on } (0, 1) \\ u(x) &= 0 && \text{for } x = 0, x = 1, \end{aligned} \tag{7.1}$$

where  $a > 0$  and  $r \geq 0$ . The basic existence and uniqueness result for general elliptic differential equations is based on Lax-Milgram's Theorem, which we will describe in section 7.3. We shall see that its stability properties, based on so called energy estimates, is automatically satisfied for finite element methods in contrast to finite difference methods.

Our goal, for a given tolerance TOL, is to find an approximation  $u_h$  of (7.1) satisfying

$$\|u - u_h\| \leq \text{TOL},$$

using few degrees of freedom by adaptive finite element approximation. Adaptive methods are based on:

- (1) an automatic mesh generator,
- (2) a numerical method ( e.g. the finite element method),

- (3) a refinement criteria (e.g. a posteriori error estimation), and
- (4) a solution algorithm ( e.g. the multigrid method).

## 7.1 The Finite Element Method

A derivation of the finite element method can be divided into:

- (1) variational formulation in an infinite dimensional space  $V$ ,
- (2) variational formulation in a finite dimensional subspace,  $V_h \subset V$ ,
- (3) choice of a basis for  $V_h$ , and
- (4) solution of the discrete system of equations.

**Step 1.** *Variational formulation in an infinite dimensional space,  $V$ .*

Consider the following Hilbert space,

$$V = \left\{ v : (0, 1) \rightarrow \mathbf{R} : \int_0^1 (v^2(x) + (v'(x))^2) dx < \infty, v(0) = v(1) = 0 \right\}.$$

Multiply equation (7.1) by  $v \in V$  and integrate by parts to get

$$\begin{aligned} \int_0^1 f v dx &= \int_0^1 ((-au')' + ru)v dx \\ &= [-au'v]_0^1 + \int_0^1 (au'v' + ruv) dx \\ &= \int_0^1 (au'v' + ruv) dx. \end{aligned} \tag{7.2}$$

Therefore the variational formulation of (7.1) is to find  $u \in V$  such that

$$A(u, v) = L(v) \quad \forall v \in V, \tag{7.3}$$

where

$$\begin{aligned} A(u, v) &= \int_0^1 (au'v' + ruv) dx, \\ L(v) &= \int_0^1 f v dx. \end{aligned}$$

**Remark 7.1.** The integration by parts in (7.2) shows that a smooth solution of equation (7.1) satisfies the variational formulation (7.3). For a solution of the variational formulation (7.3) to also be a solution of the equation (7.1), we need additional conditions



on the regularity of the functions  $a, r$  and  $f$  so that  $u''$  is continuous. Then the following integration by parts yields, as in (7.2),

$$0 = \int_0^1 (au'v' + ruv - fv) dx = \int_0^1 (-(au')' + ru - f)v dx.$$

Since this holds for all  $v \in V$ , it implies that

$$-(au')' + ru - f = 0,$$

provided  $-(au')' + ru - f$  is continuous.

**Step 2.** *Variational formulation in the finite dimensional subspace,  $V_h$ .*

First divide the interval  $(0, 1)$  into  $0 = x_0 < x_1 < \dots < x_{N+1} = 1$ , i.e. generate the mesh. Then define the space of continuous piecewise linear functions on the mesh with zero boundary conditions

$$V_h = \{v \in V \quad : \quad v(x) |_{(x_i, x_{i+1})} = c_i x + d_i, \text{ i.e. } v \text{ is linear on } (x_i, x_{i+1}), i = 0, \dots, N \\ \text{and } v \text{ is continuous on } (0, 1)\}.$$

The variational formulation in the finite dimensional subspace is to find  $u_h \in V_h$  such that

$$A(u_h, v) = L(v) \quad \forall v \in V_h. \tag{7.4}$$

The function  $u_h$  is a finite element solution of the equation (7.1). Other finite element solutions are obtained from alternative finite dimensional subspaces, e.g. based on piecewise quadratic approximation.

**Step 3.** *Choose a basis for  $V_h$ .*

Let us introduce the basis functions  $\phi_i \in V_h$ , for  $i = 1, \dots, N$ , defined by

$$\phi_i(x_j) = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j. \end{cases} \tag{7.5}$$

A function  $v \in V_h$  has the representation

$$v(x) = \sum_{i=1}^N v_i \phi_i(x),$$

where  $v_i = v(x_i)$ , i.e. each  $v \in V_h$  can be written in a unique way as a linear combination of the basis functions  $\phi_i$ .

**Step 4.** *Solve the discrete problem (7.4).*

Using the basis functions  $\phi_i$ , for  $i = 1, \dots, N$  from Step 3, we have

$$u_h(x) = \sum_{i=1}^N \xi_i \phi_i(x),$$

where  $\xi = (\xi_1, \dots, \xi_N)^T \in \mathbf{R}^N$ , and choosing  $v = \phi_j$  in (7.4), we obtain

$$\begin{aligned} L(\phi_j) &= A(u_h, \phi_j) \\ &= A\left(\sum_i \phi_i \xi_i, \phi_j\right) = \sum_i \xi_i A(\phi_i, \phi_j), \end{aligned}$$

so that  $\xi \in \mathbf{R}^N$  solves the linear system

$$\tilde{A}\xi = \tilde{L}, \quad (7.6)$$

where

$$\begin{aligned} \tilde{A}_{ji} &= A(\phi_i, \phi_j), \\ \tilde{L}_j &= L(\phi_j). \end{aligned}$$

The  $N \times N$  matrix  $\tilde{A}$  is called the stiffness matrix and the vector  $\tilde{L} \in \mathbf{R}^N$  is called the load vector.

**Example 7.2.** Consider the following two dimensional problem,

$$\begin{aligned} -\operatorname{div}(k\nabla u) + ru &= f \quad \text{in } \Omega \subset \mathbb{R}^2 \\ u &= g_1 \quad \text{on } \Gamma_1 \\ \frac{\partial u}{\partial n} &= g_2 \quad \text{on } \Gamma_2, \end{aligned} \quad (7.7)$$

where  $\partial\Omega = \Gamma = \Gamma_1 \cup \Gamma_2$  and  $\Gamma_1 \cap \Gamma_2 = \emptyset$ . The variational formulation has the following form.

(i) Variational formulation in the infinite dimensional space.

Let

$$V_g = \left\{ v(x) : \int_{\Omega} (v^2(x) + |\nabla v(x)|^2) dx < \infty, v|_{\Gamma_1} = g \right\}.$$

Take a function  $v \in V_0$ , i.e.  $v = 0$  on  $\Gamma_1$ , then by (7.7)

$$\begin{aligned} \int_{\Omega} f v dx &= - \int_{\Omega} \operatorname{div}(k\nabla u) v dx + \int_{\Omega} r u v dx \\ &= \int_{\Omega} k \nabla u \cdot \nabla v dx - \int_{\Gamma_1} k \frac{\partial u}{\partial n} v ds - \int_{\Gamma_2} k \frac{\partial u}{\partial n} v ds + \int_{\Omega} r u v dx \\ &= \int_{\Omega} k \nabla u \cdot \nabla v dx - \int_{\Gamma_2} k g_2 v ds + \int_{\Omega} r u v dx. \end{aligned}$$

The variational formulation for the model problem (7.7) is to find  $u \in V_{g_1}$  such that

$$A(u, v) = L(v) \quad \forall v \in V_0, \quad (7.8)$$

where

$$\begin{aligned} A(u, v) &= \int_{\Omega} (k \nabla u \cdot \nabla v + ruv) \, dx, \\ L(v) &= \int_{\Omega} f v \, dx + \int_{\Gamma_2} k g_2 v \, ds. \end{aligned}$$

(ii) Variational formulation in the finite dimensional space.

Assume for simplicity that  $\Omega$  is a polygonal domain which can be divided into a triangular mesh  $T_h = \{K_1, \dots, K_N\}$  of non overlapping triangles  $K_i$  and let  $h = \max_i(\text{length of longest side of } K_i)$ . Assume also that the boundary function  $g_1$  is continuous and that its restriction to each edge  $K_i \cap \Gamma_1$  is a linear function. Define

$$\begin{aligned} V_0^h &= \{v \in V_0 : v|_{K_i} \text{ is linear } \forall K_i \in T_h, v \text{ is continuous on } \Omega\}, \\ V_{g_1}^h &= \{v \in V_{g_1} : v|_{K_i} \text{ is linear } \forall K_i \in T_h, v \text{ is continuous on } \Omega\}, \end{aligned}$$

and the finite element method is to find  $u_h \in V_{g_1}^h$  such that

$$A(u_h, v) = L(v), \quad \forall v \in V_0^h. \quad (7.9)$$

(iii) Choose a basis for  $V_0^h$ .

As in the one dimensional problem, choose the basis  $\phi_j \in V_0^h$  such that

$$\phi_j(x_i) = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases} \quad j = 1, 2, \dots, N,$$

where  $x_i, i = 1, \dots, N$ , are the vertices of the triangulation.

(iv) Solve the discrete system.

Let

$$u_h(x) = \sum_{i=1}^N \xi_i \phi_i(x), \quad \text{and } \xi_i = u_h(x_i).$$

Then (7.9) can be written in matrix form,

$$\tilde{A} \xi = \tilde{L}, \quad \text{where } \tilde{A}_{ji} = A(\phi_i, \phi_j) \text{ and } \tilde{L}_j = L(\phi_j).$$

□

## 7.2 Error Estimates and Adaptivity

We shall now study a priori and a posteriori error estimates for finite element methods, where

$$\begin{aligned} \|u - u_h\| &\leq E_1(h, u, f) \quad \text{is an a priori error estimate,} \\ \|u - u_h\| &\leq E_2(h, u_h, f) \quad \text{is an a posteriori error estimate.} \end{aligned}$$

Before we start, let us study the following theorem, which we will prove later,

**Theorem 7.3** (Lax-Milgram). *Let  $V$  be a Hilbert space with norm  $\|\cdot\|_V$  and scalar product  $(\cdot, \cdot)_V$  and assume that  $A$  is a bilinear functional and  $L$  is a linear functional that satisfy:*

- (1)  $A$  is symmetric, i.e.  $A(v, w) = A(w, v) \quad \forall v, w \in V$ ;
- (2)  $A$  is  $V$ -elliptic, i.e.  $\exists \alpha > 0$  such that  $A(v, v) \geq \alpha \|v\|_V^2 \quad \forall v \in V$ ;
- (3)  $A$  is continuous, i.e.  $\exists C \in \mathbb{R}$  such that  $|A(v, w)| \leq C \|v\|_V \|w\|_V$ ; and
- (4)  $L$  is continuous, i.e.  $\exists \Lambda \in \mathbb{R}$  such that  $|L(v)| \leq \Lambda \|v\|_V \quad \forall v \in V$ .

Then there is a unique function  $u \in V$  such that  $A(u, v) = L(v) \quad \forall v \in V$ , and the stability estimate  $\|u\|_V \leq \Lambda/\alpha$  holds.

### 7.2.1 An A Priori Error Estimate

The approximation property of the space  $V_h$  can be characterized by

**Lemma 7.4.** *Suppose  $V_h$  is the piecewise linear finite element space (7.4), which discretizes the functions in  $V$ , defined on  $(0, 1)$ , with the interpolant  $\pi : V \rightarrow V_h$  defined by*

$$\pi v(x) = \sum_{i=1}^N v(x_i) \phi_i(x), \quad (7.10)$$

where  $\{\phi_i\}$  is the basis (7.5) of  $V_h$ . Then

$$\begin{aligned} \|(v - \pi v)'\|_{L^2(0,1)} &\leq \sqrt{\int_0^1 h^2 v''(x)^2 dx} \leq Ch, \\ \|v - \pi v\|_{L^2(0,1)} &\leq \sqrt{\int_0^1 h^4 v''(x)^2 dx} \leq Ch^2, \end{aligned} \quad (7.11)$$

where  $h = \max_i (x_{i+1} - x_i)$ .

*Proof.* Take  $v \in V$  and consider first (7.11) on an interval  $(x_i, x_{i+1})$ . By the mean value theorem, there is for each  $x \in (x_i, x_{i+1})$  a  $\xi \in (x_i, x_{i+1})$  such that  $v'(\xi) = (\pi v)'(x)$ . Therefore

$$v'(x) - (\pi v)'(x) = v'(x) - v'(\xi) = \int_{\xi}^x v''(s) ds,$$

so that

$$\begin{aligned} \int_{x_i}^{x_{i+1}} |v'(x) - (\pi v)'(x)|^2 dx &= \int_{x_i}^{x_{i+1}} \left( \int_{\xi}^x v''(s) ds \right)^2 dx \\ &\leq \int_{x_i}^{x_{i+1}} |x - \xi| \int_{\xi}^x (v''(s))^2 ds dx \\ &\leq h^2 \int_{x_i}^{x_{i+1}} (v''(s))^2 ds, \end{aligned} \quad (7.12)$$

which after summation of the intervals proves (7.11).

Next, we have

$$v(x) - \pi v(x) = \int_{x_i}^x (v - \pi v)'(s) ds,$$

so by (7.12)

$$\begin{aligned} \int_{x_i}^{x_{i+1}} |v(x) - \pi v(x)|^2 dx &= \int_{x_i}^{x_{i+1}} \left( \int_{x_i}^x (v - \pi v)'(s) ds \right)^2 dx \\ &\leq \int_{x_i}^{x_{i+1}} |x - x_i| \int_{x_i}^x ((v - \pi v)')^2(s) ds dx \\ &\leq h^4 \int_{x_i}^{x_{i+1}} (v''(s))^2 ds, \end{aligned}$$

which after summation of the intervals proves the lemma.  $\square$

Our derivation of the a priori error estimate

$$\|u - u_h\|_V \leq Ch,$$

where  $u$  and  $u_h$  satisfy (7.3) and (7.4), respectively, uses Lemma 7.4 and a combination of the following four steps:

(1) error representation based on the *ellipticity*

$$\alpha \int_{\Omega} (v^2(x) + (v'(x))^2) dx \leq A(v, v) = \int_{\Omega} (a(v')^2 + rv^2) dx,$$

where  $\alpha = \inf_{x \in (0,1)} (a(x), r(x)) > 0$ ,

(2) the *orthogonality*

$$A(u - u_h, v) = 0 \quad \forall v \in V_h,$$

obtained by  $V_h \subset V$  and subtraction of the two equations

$$\begin{aligned} A(u, v) &= L(v) \quad \forall v \in V \quad \text{by (7.3),} \\ A(u_h, v) &= L(v) \quad \forall v \in V_h \quad \text{by (7.4),} \end{aligned}$$

(3) the *continuity*

$$|A(v, w)| \leq C \|v\|_V \|w\|_V \quad \forall v, w \in V,$$

where  $C \leq \sup_{x \in (0,1)} (a(x), r(x))$ , and

(4) the *interpolation estimates*

$$\begin{aligned} \|(v - \pi v)'\|_{L^2} &\leq Ch, \\ \|v - \pi v\|_{L^2} &\leq Ch^2, \end{aligned} \tag{7.13}$$

where  $h = \max (x_{i+1} - x_i)$ .

To start the proof of an a priori estimate let  $e \equiv u - u_h$ . Then by Cauchy's inequality

$$\begin{aligned} A(e, e) &= A(e, u - \pi u + \pi u - u_h) \\ &= A(e, u - \pi u) + A(e, \pi u - u_h) \\ &\stackrel{\text{Step2}}{=} A(e, u - \pi u) \\ &\leq \sqrt{A(e, e)} \sqrt{A(u - \pi u, u - \pi u)}, \end{aligned}$$

so that by division of  $\sqrt{A(e, e)}$ ,

$$\begin{aligned} \sqrt{A(e, e)} &\leq \sqrt{A(u - \pi u, u - \pi u)} \\ &\stackrel{\text{Step3}}{=} C \|u - \pi u\|_V \\ &\equiv C \sqrt{\|u - \pi u\|_{L^2}^2 + \|(u - \pi u)'\|_{L^2}^2} \\ &\stackrel{\text{Step4}}{\leq} Ch. \end{aligned}$$

Therefore, by Step 1

$$\alpha \|e\|_V^2 \leq A(e, e) \leq Ch^2,$$

which implies the a priori estimate

$$\|e\|_V \leq Ch,$$

where  $C = K(u)$ . □

### 7.2.2 An A Posteriori Error Estimate

**Example 7.5.** Consider the model problem (7.1), namely,

$$\begin{cases} -(au')' + ru = f & \text{in } (0, 1), \\ u(0) = u(1) = 0. \end{cases}$$

Then

$$\begin{aligned} \sqrt{A(u - u_h, u - u_h)} &\leq C \|a^{-\frac{1}{2}}(f - ru_h + a'u_h')h\|_{L^2} \\ &\equiv E(h, u_h, f). \end{aligned} \tag{7.14}$$

*Proof.* Let  $e = u - u_h$  and let  $\pi e \in V_h$  be the nodal interpolant of  $e$ . We have

$$\begin{aligned} A(e, e) &= A(e, e - \pi e) \quad (\text{by orthogonality}) \\ &= A(u, e - \pi e) - A(u_h, e - \pi e). \end{aligned}$$

Using the notation  $(f, v) \equiv \int_0^1 f v \, dx$ , we obtain by integration by parts

$$\begin{aligned} A(e, e) &= (f, e - \pi e) - \sum_{i=1}^N \int_{x_i}^{x_{i+1}} (au'_h(e - \pi e)' + ru_h(e - \pi e)) \, dx \\ &= (f - ru_h, e - \pi e) - \sum_{i=1}^N \left\{ [au'_h(e - \pi e)]_{x_i}^{x_{i+1}} - \int_{x_i}^{x_{i+1}} (au'_h)'(e - \pi e) \, dx \right\} \\ &= (f - ru_h + a'u'_h, e - \pi e) \quad (\text{since } u_h''|_{(x_i, x_{i+1})} = 0, (e - \pi e)(x_i) = 0) \\ &\leq \|a^{-\frac{1}{2}}h(f - ru_h + a'u'_h)\|_{L^2} \|a^{\frac{1}{2}}h^{-1}(e - \pi e)\|_{L^2}. \end{aligned}$$

Lemma 7.6 implies

$$\sqrt{A(e, e)} \leq C \|a^{-\frac{1}{2}}h(f - ru_h + a'u'_h)\|_{L^2},$$

which also shows that

$$\|e\|_V \leq Ch,$$

where  $C = K'(u_h)$ . □

**Lemma 7.6.** *There is a constant  $C$ , independent of  $u$  and  $u_h$ , such that,*

$$\|a^{\frac{1}{2}}h^{-1}(e - \pi e)\|_{L^2} \leq C \sqrt{\int_0^1 ae'e' \, dx} \leq C \sqrt{A(e, e)}$$

**Exercise 7.7.** Use the interpolation estimates in Lemma 7.4 to prove Lemma 7.6.

### 7.2.3 An Adaptive Algorithm

We formulate an adaptive algorithm based on the a posteriori error estimate (7.14) as follows:

- (1) Choose an initial coarse mesh  $T_{h_0}$  with mesh size  $h_0$ .
- (2) Compute the corresponding FEM solution  $u_{h_i}$  in  $V_{h_i}$ .
- (3) Given a computed solution  $u_{h_i}$  in  $V_{h_i}$ , with the mesh size  $h_i$ ,

stop                    if  $E(h_i, u_{h_i}, f) \leq TOL$   
go to step 4        if  $E(h_i, u_{h_i}, f) > TOL$ .

(4) Determine a new mesh  $T_{h_{i+1}}$  with mesh size  $h_{i+1}$  such that

$$E(h_{i+1}, u_{h_i}, f) \cong TOL,$$

by letting the error contribution for all elements be approximately constant, i.e.

$$\|a^{-\frac{1}{2}}h(f - ru_h - a'u'_h)\|_{L^2(x_i, x_{i+1})} \cong C, \quad i = 1, \dots, N,$$

then go to Step 2.

### 7.3 Lax-Milgram's Theorem

**Theorem 7.8.** *Suppose  $A$  is symmetric, i.e.  $A(u, v) = A(v, u) \quad \forall u, v \in V$ , then (Variational problem)  $\iff$  (Minimization problem) with*

$$\begin{aligned} (\text{Var}) \quad & \text{Find } u \in V \text{ such that } A(u, v) = L(v) \quad \forall v \in V, \\ (\text{Min}) \quad & \text{Find } u \in V \text{ such that } F(u) \leq F(v) \quad \forall v \in V, \end{aligned}$$

where

$$F(w) \equiv \frac{1}{2}A(w, w) - L(w) \quad \forall w \in V.$$

*Proof.* Take  $\epsilon \in \mathbb{R}$ . Then

$$\begin{aligned} (\Rightarrow) \quad F(u + \epsilon w) &= \frac{1}{2}A(u + \epsilon w, u + \epsilon w) - L(u + \epsilon w) \\ &= \left( \frac{1}{2}A(u, u) - L(u) \right) + \epsilon A(u, w) - \epsilon L(w) + \frac{1}{2}\epsilon^2 A(w, w) \\ &\geq \left( \frac{1}{2}A(u, u) - L(u) \right) \quad \left( \text{since } \frac{1}{2}\epsilon^2 A(w, w) \geq 0 \text{ and } A(u, w) = L(w) \right) \\ &= F(u). \end{aligned}$$

( $\Leftarrow$ ) Let  $g(\epsilon) = F(u + \epsilon w)$ , where  $g : \mathbf{R} \rightarrow \mathbf{R}$ . Then

$$0 = g'(0) = 0 \cdot A(w, w) + A(u, w) - L(w) = A(u, w) - L(w).$$

Therefore

$$A(u, w) = L(w) \quad \forall w \in V.$$

□

**Theorem 7.9** (Lax-Milgram). *Let  $V$  be a Hilbert space with norm  $\|\cdot\|_V$  and scalar product  $(\cdot, \cdot)_V$  and assume that  $A$  is a bilinear functional and  $L$  is a linear functional that satisfy:*

- (1)  $A$  is symmetric, i.e.  $A(v, w) = A(w, v) \quad \forall v, w \in V$ ;
- (2)  $A$  is  $V$ -elliptic, i.e.  $\exists \alpha > 0$  such that  $A(v, v) \geq \alpha \|v\|_V^2 \quad \forall v \in V$ ;



(3)  $A$  is continuous, i.e.  $\exists C \in \mathbb{R}$  such that  $|A(v, w)| \leq C\|v\|_V\|w\|_V$ ; and

(4)  $L$  is continuous, i.e.  $\exists \Lambda \in \mathbb{R}$  such that  $|L(v)| \leq \Lambda\|v\|_V \quad \forall v \in V$ .

Then there is a unique function  $u \in V$  such that  $A(u, v) = L(v) \quad \forall v \in V$ , and the stability estimate  $\|u\|_V \leq \Lambda/\alpha$  holds.

*Proof.* The goal is to construct  $u \in V$  solving the minimization problem  $F(u) \leq F(v)$  for all  $v \in V$ , which by the previous theorem is equivalent to the variational problem. The energy norm,  $\|v\|^2 \equiv A(v, v)$ , is equivalent to the norm of  $V$ , since by Condition 2 and 3,

$$\alpha\|v\|_V^2 \leq A(v, v) = \|v\|^2 \leq C\|v\|_V^2.$$

Let

$$\beta = \inf_{v \in V} F(v). \quad (7.15)$$

Then  $\beta \in \mathbf{R}$ , since

$$F(v) = \frac{1}{2}\|v\|^2 - L(v) \geq \frac{1}{2}\|v\|^2 - \Lambda\|v\| \geq -\frac{\Lambda^2}{2}.$$

We want to find a solution to the minimization problem  $\min_{v \in V} F(v)$ . It is therefore natural to study a minimizing sequence  $v_i$ , such that

$$F(v_i) \rightarrow \beta = \inf_{v \in V} F(v). \quad (7.16)$$

The next step is to conclude that the  $v_i$  infact converge to a limit:

$$\begin{aligned} \left\| \frac{v_i - v_j}{2} \right\|^2 &= \frac{1}{2}\|v_i\|^2 + \frac{1}{2}\|v_j\|^2 - \left\| \frac{v_i + v_j}{2} \right\|^2 \quad (\text{by the parallelogram law}) \\ &= \frac{1}{2}\|v_i\|^2 - L(v_i) + \frac{1}{2}\|v_j\|^2 - L(v_j) \\ &\quad - \left( \left\| \frac{v_i + v_j}{2} \right\|^2 - 2L\left(\frac{v_i + v_j}{2}\right) \right) \\ &= F(v_i) + F(v_j) - 2F\left(\frac{v_i + v_j}{2}\right) \\ &\leq F(v_i) + F(v_j) - 2\beta \quad (\text{by (7.15)}) \\ &\rightarrow 0, \quad (\text{by (7.16)}). \end{aligned}$$

Hence  $\{v_i\}$  is a Cauchy sequence in  $V$  and since  $V$  is a Hilbert space ( in particular  $V$  is a complete space) we have  $v_i \rightarrow u \in V$ .

Finally  $F(u) = \beta$ , since

$$\begin{aligned}
 |F(v_i) - F(u)| &= \left| \frac{1}{2}(\|v_i\|^2 - \|u\|^2) - L(v_i - u) \right| \\
 &= \left| \frac{1}{2}A(v_i - u, v_i + u) - L(v_i - u) \right| \\
 &\leq \left( \frac{C}{2}\|v_i + u\|_V + \Lambda \right) \|v_i - u\|_V \\
 &\rightarrow 0.
 \end{aligned}$$

Therefore there exists a unique (why?) function  $u \in V$  such that  $F(u) \leq F(v) \quad \forall v \in V$ . To verify the stability estimate, take  $v = u$  in (Var) and use the ellipticity (1) and continuity (3) to obtain

$$\alpha \|u\|_V^2 \leq A(u, u) = L(u) \leq \Lambda \|u\|_V$$

so that

$$\|u\|_V \leq \frac{\Lambda}{\alpha}.$$

The uniqueness of  $u$  can also be verified from the stability estimate. If  $u_1, u_2$  are two solutions of the variational problem we have  $A(u_1 - u_2, v) = 0$  for all  $v \in V$ . Therefore the stability estimate implies  $\|u_1 - u_2\|_V = 0$ , i.e.  $u_1 = u_2$  and consequently the solution is unique.  $\square$

**Example 7.10.** Determine conditions for the functions  $k, r$  and  $f : \Omega \rightarrow \mathbb{R}$  such that the assumptions in the Lax-Milgram theorem are satisfied for the following elliptic partial differential equation in  $\Omega \subset \mathbf{R}^2$

$$\begin{aligned}
 -\operatorname{div}(k\nabla u) + ru &= f \quad \text{in } \Omega \\
 u &= 0 \quad \text{on } \partial\Omega.
 \end{aligned}$$

**Solution.** This problem satisfies (Var) with

$$V = \{v : \int_{\Omega} (v^2(x) + |\nabla v(x)|^2) dx < \infty, \text{ and } v|_{\partial\Omega} = 0\},$$

$$\begin{aligned}
 A(u, v) &= \int_{\Omega} (k\nabla u \nabla v + ruv) dx, \\
 L(v) &= \int_{\Omega} fv dx, \\
 \|v\|_V^2 &= \int_{\Omega} (v^2(x) + |\nabla v|^2) dx.
 \end{aligned}$$

Consequently  $V$  is a Hilbert space and  $A$  is symmetric and continuous provided  $k$  and  $r$  are uniformly bounded.

The ellipticity follows by

$$\begin{aligned} A(v, v) &= \int_{\Omega} (k|\nabla v|^2 + rv^2) dx \\ &\geq \alpha \int_{\Omega} (v^2(x) + |\nabla v|^2) dx \\ &= \alpha \|v\|_{H^1}^2, \end{aligned}$$

provided  $\alpha = \inf_{x \in \Omega} (k(x), r(x)) > 0$ .

The continuity of  $A$  is a consequence of

$$\begin{aligned} A(v, w) &\leq \max(\|k\|_{L^\infty}, \|r\|_{L^\infty}) \int_{\Omega} (|\nabla v| |\nabla w| + |v| |w|) dx \\ &\leq \max(\|k\|_{L^\infty}, \|r\|_{L^\infty}) \|v\|_{H^1} \|w\|_{H^1}, \end{aligned}$$

provided  $\max(\|k\|_{L^\infty}, \|r\|_{L^\infty}) = C < \infty$ .

Finally, the functional  $L$  is continuous, since

$$|L(v)| \leq \|f\|_{L^2} \|v\|_{L^2} \leq \|f\|_{L^2} \|v\|_V,$$

which means that we may take  $\Lambda = \|f\|_{L^2}$  provided we assume that  $f \in L^2(\Omega)$ . Therefore the problem satisfies the Lax-Milgram theorem.  $\square$

**Example 7.11.** Verify that the assumption of the Lax-Milgram theorem are satisfied for the following problem,

$$\begin{aligned} -\Delta u &= f \quad \text{in } \Omega, \\ u &= 0 \quad \text{on } \partial\Omega. \end{aligned}$$

**Solution.** This problem satisfies (Var) with

$$\begin{aligned} V = H_0^1 &= \{v \in H^1 : v|_{\partial\Omega} = 0\}, \\ H^1 &= \{v : \int_{\Omega} (v^2(x) + |\nabla v(x)|^2) dx < \infty\}, \end{aligned}$$

$$\begin{aligned} A(u, v) &= \int_{\Omega} \nabla u \nabla v dx, \\ L(v) &= \int_{\Omega} f v dx. \end{aligned}$$

To verify the V-ellipticity, we use the *Poincaré inequality*, i.e. there is a constant  $C$  such that

$$v \in H_0^1 \Rightarrow \int_{\Omega} v^2 dx \leq C \int_{\Omega} |\nabla v|^2 dx. \quad (7.17)$$

In one dimension and  $\Omega = (0, 1)$ , the inequality (7.17) takes the form

$$\int_0^1 v^2(x) \, dx \leq \int_0^1 (v'(x))^2 \, dx, \quad (7.18)$$

provided  $v(0) = 0$ . Since

$$v(x) = v(0) + \int_0^x v'(s) \, ds = \int_0^x v'(s) \, ds,$$

and by Cauchy's inequality

$$\begin{aligned} v^2(x) &= \left( \int_0^x v'(s) \, ds \right)^2 \leq x \int_0^x v'(s)^2 \, ds \\ &\leq \int_0^1 v'(s)^2 \, ds \quad \text{since } x \in (0, 1). \end{aligned}$$

The V-ellipticity of A follows by (7.18) and

$$\begin{aligned} A(v, v) &= \int_0^1 v'(x)^2 \, dx = \frac{1}{2} \int_0^1 \left( (v'(x))^2 \, dx + \frac{1}{2}(v'(x))^2 \right) \, dx \\ &\geq \frac{1}{2} \int_0^1 (v'(x)^2 + v(x)^2) \, dx \\ &= \frac{1}{2} \|v\|_{H_0^1}^2 \quad \forall v \in H_0^1. \end{aligned}$$

The other conditions can be proved similarly as in the previous example. Therefore this problem satisfies the Lax-Milgram theorem.  $\square$

## Chapter 8

# Markov Chains, Duality and Dynamic Programming, by Jonathan Goodman

### 8.1 Introduction

There are two main ideas in the arbitrage theory of pricing. One is that in complete markets, everyone should agree on a common price – any other price leads to an arbitrage opportunity. The other is that this price is the expected value of the cash flow with respect to some probability model – risk neutral pricing. In the simplest case, this probability model is a discrete Markov chain. This lecture describes how to compute probabilities and expected values for discrete Markov chain models. This is the main computational step in "risk neutral" option pricing.

The methods here compute the expected values by a time marching process that uses the transition matrix. Another evolution process allows us to compute probabilities. These evolution processes are related but not the same. The relation between the forward evolution for probabilities and the backward evolution for expected values is called *duality*. It is similar to the relation between a matrix and its transpose. The transpose of a matrix is sometimes called its dual.

The method of risk neutral arbitrage pricing extends to other more technical situations, but the main ideas are clear in the simple context of Markov chains. If the Markov chain model is replaced by a stochastic differential equation model, then the transition matrix is replaced by a partial differential operator – the "generator", and the matrix transpose is replaced by the "dual" of this generator. This is the subject of future lectures.

Many financial instruments allow the holder to make decisions along the way that effect the ultimate value of the instrument. American style options, loans that be repaid early, and convertible bonds are examples. To compute the value of such an instrument, we also seek the optimal decision strategy. *Dynamic programming* is a computational method that computes the value and decision strategy at the same time. It reduces the difficult "multiperiod decision problem" to a sequence of hopefully easier "single

period” problems. It works backwards in time much as the expectation method does. The tree method commonly used to value American style stock options is an example of the general dynamic programming method.

## 8.2 Markov Chains

(This section assumes familiarity with basic probability theory using mathematicians’ terminology. References on this include the probability books by G. C. Rota, W. Feller, Hoel and Stone, and B. V. Gnedenko.)

Many discrete time discrete state space stochastic models are stationary discrete Markov chains. Such a Markov chain is characterized by its state space,  $\mathcal{S}$ , and its transition matrix,  $P$ . We use the following notations:

- $x, y, \dots$ : possible states of the system, elements of  $\mathcal{S}$ .
- The possible times are  $t = 0, 1, 2, \dots$
- $X(t)$ : the (unknown) state of the system at time  $t$ . It is some element of  $\mathcal{S}$ .
- $u(x, t) = \mathbf{Pr}(X(t) = x)$ . These probabilities satisfy an evolution equation moving forward in time. We use similar notation for conditional probabilities, for example,  $u(x, t|X(0) = x_0) = \mathbf{Pr}(X(t) = x|X(0) = x_0)$ .
- $p(x, y) = \mathbf{Pr}(x \rightarrow y) = \mathbf{Pr}(X(t+1) = y|X(t) = x)$ . These “transition probabilities” are the elements of the transition matrix,  $P$ .

The transition probabilities have the properties:

$$0 \leq p(x, y) \leq 1 \quad \text{for all } x \in \mathcal{S} \text{ and } y \in \mathcal{S}. \quad (8.1)$$

and

$$\sum_{y \in \mathcal{S}} p(x, y) = 1 \quad \text{for all } x \in \mathcal{S}. \quad (8.2)$$

The first is because the  $p(x, y)$  are probabilities, the second because the state  $x$  must go somewhere, possibly back to  $x$ . It is not true that

$$\text{(NOT ALWAYS TRUE)} \quad \sum_{x \in \mathcal{S}} p(x, y) = 1 \quad . \quad \text{(NOT ALWAYS TRUE)}$$

The Markov property is that knowledge of the state at time  $t$  is all the information about the present and past relevant to predicting the future. That is:

$$\begin{aligned} \mathbf{Pr}(X(t+1) = y|X(t) = x_0, X(t-1) = x_1, \dots) \\ = \mathbf{Pr}(X(t+1) = y|X(t) = x_0) \end{aligned} \quad (8.3)$$

no matter what extra history information ( $X(t-1) = x_1, \dots$ ) we have. This may be thought of as a lack of long term memory. It may also be thought of as a completeness

property of the model: the state space is rich enough to characterize the state of the system at time  $t$  completely.

To illustrate this point, consider the model

$$Z(t+1) = aZ(t) + bZ(t-1) + \xi(t), \quad (8.4)$$

where the  $\xi(t)$  are independent random variables. Models like this are used in “time series analysis”. Here  $Z$  is a continuous variable instead a discrete variable to make the example simpler. If we say that the state at time  $t$  is  $Z(t)$  then (8.4) is not a Markov chain. Clearly we do better at predicting  $Z(t+1)$  if we know both  $Z(t)$  and  $Z(t-1)$  than if we know just  $Z(t)$ . If we say that the state at time  $t$  is the two dimensional vector

$$X(t) = \begin{pmatrix} Z(t) \\ Z(t-1) \end{pmatrix},$$

then

$$\begin{pmatrix} Z(t) \\ Z(t-1) \end{pmatrix} = \begin{pmatrix} a & b \\ 1 & 0 \end{pmatrix} \begin{pmatrix} Z(t-1) \\ Z(t-2) \end{pmatrix} + \begin{pmatrix} \xi(t) \\ 0 \end{pmatrix}$$

may be rewritten

$$X(t+1) = AX(t) + \begin{pmatrix} \xi(t) \\ 0 \end{pmatrix}.$$

Thus,  $X(t)$  is a Markov chain. This trick of expressing lag models with multidimensional states is common in time series analysis.

The simpler of the evolutions, and the one less used in practice, is the forward evolution for the probabilities  $u(x, t)$ . Once we know the numbers  $u(x, t)$  for all  $x \in \mathcal{S}$  and a particular  $t$ , we can compute them for  $t+1$ . Proceeding in this way, starting from the numbers  $u(x, 0)$  for all  $x \in \mathcal{S}$ , we can compute up to whatever  $T$  is desired. The evolution equation for the probabilities  $u(x, t)$  is found using conditional probability:

$$\begin{aligned} u(x, t+1) &= \mathbf{Pr}(X(t+1) = x) \\ &= \sum_{y \in \mathcal{S}} \mathbf{Pr}(X(t+1) = x | X(t) = y) \cdot \mathbf{Pr}(X(t) = y) \\ u(x, t+1) &= \sum_{y \in \mathcal{S}} p(y, x) u(y, t). \end{aligned} \quad (8.5)$$

To express this in matrix form, we suppose that the state space,  $\mathcal{S}$ , is finite, and that the states have been numbered  $x_1, \dots, x_n$ . The transition matrix,  $P$ , is  $n \times n$  and has  $(i, j)$  entry  $p_{ij} = p(x_i, x_j)$ . We sometimes conflate  $i$  with  $x_i$  and write  $p_{xy} = p(x, y)$ ; until you start programming the computer, there is no need to order the states. With this convention, (8.5) can be interpreted as vector–matrix multiplication if we define a row vector  $\underline{u}(t)$  with components  $(u_1(t), \dots, u_n(t))$ , where we have written  $u_i(t)$  for  $u(x_i, t)$ . As long as ordering is unimportant, we could also write  $u_x(t) = u(x, t)$ . Now, (8.5) can be rewritten

$$\underline{u}(t+1) = \underline{u}(t)P. \quad (8.6)$$

Since  $\underline{u}$  is a row vector, the expression  $P\underline{u}$  does not make sense because the dimensions of the matrices are incompatible for matrix multiplication. The convention of using a row vector for the probabilities and therefore putting the vector in the left of the matrix is common in applied probability. The relation (8.6) can be used repeatedly<sup>1</sup>

$$\begin{aligned} \underline{u}(1) &= \underline{u}(0)P \text{ and } \underline{u}(2) = \underline{u}(1)P \\ &\quad \rightarrow \\ \underline{u}(2) &= (\underline{u}(0)P)P = \underline{u}(0)(PP) = \underline{u}(0)P^2 \end{aligned}$$

to yield

$$\underline{u}(t) = \underline{u}(0)P^t, \quad (8.7)$$

where  $P^t$  means  $P$  to the power  $t$ , not the transpose of  $P$ .

Actually, the Markov property is a bit stronger than (8.3). It applies not only to events determined by time  $t + 1$ , but to any events determined in the future of  $t$ . For example, if  $A$  is the event  $X(t + 3) = x$  or  $y$  and  $X(t + 1) \neq X(t + 4)$ , then

$$\Pr(A \mid X(t) = z \text{ and } X(t - 1) = w) = \Pr(A \mid X(t) = z).$$

### 8.3 Expected Values

The more general and useful evolution equation is the backward evolution for expected values. In the simplest situation, suppose that  $X(t)$  is a Markov chain, that the probability distribution  $u(x, 0) = \Pr(X(0) = x)$  is known, and that we want to evaluate  $\mathbf{E}(V(X(T)))$ . We will call time  $t = 0$  the present, time  $t = T$  the payout time, and times  $t = 1, \dots, T - 1$  intermediate times.

The backward evolution computed the desired expected value in terms of a collection of other conditional expected values,  $f(x, t)$ , where  $x \in \mathcal{S}$  and  $t$  is an intermediate time. We start with the final time values  $f(x, T) = V(x)$  for all  $x \in \mathcal{S}$ . We then compute the numbers  $f(x, T - 1)$  using the  $f(x, t)$  and  $P$ . We continue in this way back to time  $t = 0$ .

The  $f(x, t)$  are expected values of the payout, given knowledge of the state at a future intermediate time:

$$f(x, t) = \mathbf{E}[V(X(T)) \mid X(t) = x]. \quad (8.8)$$

Recall our convention that time 0 is the present time, time  $t > 0$  is in the future, but not as far in the future as the time,  $T$ , at which the payout is made. We may think of the  $f(x, t)$  as possible expected values at the future intermediate time  $t$ . At time  $t$  we would know the value of  $X(t)$ . If that value were  $x$ , then the expected value of  $V(X(T))$  would be  $f(x, t)$ .

Instead of computing  $f(x, t)$  directly from the definition (8.8), we can compute it in terms of the  $f(x, t + 1)$  using the transition matrix. If the system is in state  $x$  at time  $t$ ,

---

<sup>1</sup>The most important fact in linear algebra is that matrix multiplication is associative:  $(AB)C = A(BC)$  for any three matrices of any size, including row or column vectors, as long as the multiplication is compatible.



then the probability for it to be at state  $y$  at the next time is  $p(x \rightarrow y) = p(x, y)$ . For expectation values, this implies

$$\begin{aligned}
 f(x, t) &= \mathbf{E}[f_T(X(T)) | X(t) = x] \\
 &= \sum_{y \in \mathcal{S}} \mathbf{E}[f_T(X(T)) | X(t+1) = y] \cdot \mathbf{Pr}(X(t+1) = y | X(t) = x) \\
 f(x, t) &= \sum_{y \in \mathcal{S}} f(y, t+1) p(x, y) .
 \end{aligned} \tag{8.9}$$

It is clear from (8.8) that  $f(x, T) = V(x)$ ; if we know the state at time  $T$  then we know the payout exactly. From these, we compute all the numbers  $f(x, T-1)$  using (8.9) with  $t = T-1$ . Continuing like this, we eventually get to  $t = 0$ . We may know  $X(0)$ , the state of the system at the current time. For example, if  $X(t)$  is the price of a stock at time  $t$ , then  $X(0) = x_0$  is the current spot price. Then the desired expected value would be  $f(x_0, 0)$ . Otherwise we can use

$$\begin{aligned}
 \mathbf{E}[V(X(T))] &= \sum_{x \in \mathcal{S}} \mathbf{E}[V(X(T)) | X(0) = x] \cdot \mathbf{Pr}(X(0) = x) \\
 &= \sum_{x \in \mathcal{S}} f(x, 0) u(x, 0) .
 \end{aligned}$$

All the values on the bottom line should be known.

Another remark on the interpretation of (8.9) will be helpful. Suppose we are at state  $x$  at time  $t$  and wish to know the expected value of  $V(X(T))$ . In one time step, starting from state  $x$ , we could go to state  $y$  at time  $t+1$  with probability<sup>2</sup>  $p(x, y)$ . The right side of (8.9) is the average over the possible  $y$  values, using probability  $p(x, y)$ . The quantities being averaged,  $f(y, t+1)$  are themselves expected values of  $V(X(T))$ . Thus, we can read (8.9) as saying that the expected value is the expected value of the expected values at the next time. A simple model for this situation is that we toss a coin. With probability  $p$  we get payout  $U$  and with probability  $1-p$  we get payout  $V$ . Let us suppose that both  $U$  and  $V$  are random with expected values  $f_U = \mathbf{E}(U)$  and  $f_V = \mathbf{E}(V)$ . The overall expected payout is  $p \cdot f_U + (1-p) \cdot f_V$ . The Markov chain situation is like this. We are at a state  $x$  at time  $t$ . We first choose state  $y \in \mathcal{S}$  with probability  $p(x, y)$ . For each  $y$  at time  $t+1$  there is a payout probability,  $U_y$ , whose probability distribution depends on  $y, t+1, V$ , and the Markov chain. The overall expected payout is the average of the expected values of the  $U_y$ , which is what (8.9) says.

As with the probability evolution equation (8.5), the equation for the evolution of the expectation values (8.9) can be written in matrix form. The difference from the probability evolution equation is that here we arrange the numbers  $f_j = f(x_j, t)$  into a *column* vector,  $\underline{f}(t)$ . The evolution equation for the expectation values is then written in matrix form as

$$\underline{f}(t) = P \underline{f}(t+1) . \tag{8.10}$$

---

<sup>2</sup>Here we should think of  $y$  as the variable and  $x$  as a parameter.

This time, the vector goes on the right. If apply (8.10) repeatedly, we get, in place of (8.7),

$$\underline{f}(t) = P^{T-t} \underline{f}(T) . \quad (8.11)$$

There are several useful variations on this theme. For example, suppose that we have a running payout rather than a final time payout. Call this payout  $g(x, t)$ . If  $X(t) = x$  then  $g(x, t)$  is added to the total payout that accumulates over time from  $t = 0$  to  $t = T$ . We want to compute

$$\mathbf{E} \left[ \sum_{t=0}^T g(X(t), t) \right] .$$

As before, we find this by computing more specific expected values:

$$f(x, t) = \mathbf{E} \left[ \sum_{t'=t}^T g(X(t'), t') | X(t) = x \right] .$$

These numbers are related through a generalization of (8.9) that takes into account the known contribution to the sum from the state at time  $t$ :

$$f(x, t) = \sum_{y \in \mathcal{S}} f(y, t+1) p(x, y) + g(x, t) .$$

The “initial condition”, given at the final time, is

$$f(x, T) = g(x, T) .$$

This includes the previous case, we take  $g(x, T) = f_T(x)$  and  $g(x, t) = 0$  for  $t < T$ .

As a final example, consider a path dependent discounting. Suppose for a state  $x$  at time  $t$  there is a discount factor  $r(x, t)$  in the range  $0 \leq r(x, t) \leq 1$ . A cash flow worth  $f$  at time  $t + 1$  will be worth  $r(x, t)f$  at time  $t$  if  $X(t) = x$ . We want the discounted value at time  $t = 0$  at state  $X(0) = x$  of a final time payout worth  $f_T(X(T))$  at time  $T$ . Define  $f(x, t)$  to be the value at time  $t$  of this payout, given that  $X(t) = x$ . If  $X(t) = x$  then the time  $t + 1$  expected discounted (to time  $t + 1$ ) value is

$$\sum_{y \in \mathcal{S}} f(y, t+1) p(x, y) .$$

This must be discounted to get the time  $t$  value, the result being

$$f(x, t) = r(x, t) \sum_{y \in \mathcal{S}} f(y, t+1) p(x, y) .$$

## 8.4 Duality and Qualitative Properties

The forward evolution equation (8.5) and the backward equation (8.9) are connected through a duality relation. For any time  $t$ , we compute (8.8) as

$$\begin{aligned} \mathbf{E} [V(X(T))] &= \sum_{x \in \mathcal{S}} \mathbf{E} [V(X(T)) | X(t) = x] \cdot \mathbf{Pr}(X(t) = x) \\ &= \sum_{x \in \mathcal{S}} f(x, t) u(x, t) . \end{aligned} \quad (8.12)$$

For now, the main point is that the sum on the bottom line does not depend on  $t$ . Given the constancy of this sum and the  $u$  evolution equation (8.5), we can give another derivation of the  $f$  evolution equation (8.9). Start with

$$\sum_{x \in \mathcal{S}} f(x, t+1)u(x, t+1) = \sum_{y \in \mathcal{S}} f(y, t)u(y, t) .$$

Then use (8.5) on the left side and rearrange the sum:

$$\sum_{y \in \mathcal{S}} \left( \sum_{x \in \mathcal{S}} f(x, t+1)p(y, x) \right) u(y, t) = \sum_{y \in \mathcal{S}} f(y, t)u(y, t) .$$

Now, if this is going to be true for any  $u(y, t)$ , the coefficients of  $u(y, t)$  on the left and right sides must be equal for each  $y$ . This gives (8.9). Similarly, it is possible to derive (8.5) from (8.9) and the constancy of the expected value.

The evolution equations (8.5) and (8.9) have some qualitative properties in common. The main one being that they preserve positivity. If  $u(x, t) \geq 0$  for all  $x \in \mathcal{S}$ , then  $u(x, t+1) \geq 0$  for all  $x \in \mathcal{S}$  also. Likewise, if  $f(x, t+1) \geq 0$  for all  $x$ , then  $f(x, t) \geq 0$  for all  $x$ . These properties are simple consequences of (8.5) and (8.9) and the positivity of the  $p(x, y)$ . Positivity preservation does not work in reverse. It is possible, for example, that  $f(x, t+1) < 0$  for some  $x$  even though  $f(x, t) \geq 0$  for all  $x$ .

The probability evolution equation (8.5) has a conservation law not shared by (8.9). It is

$$\sum_{x \in \mathcal{S}} u(x, t) = \text{const} . \quad (8.13)$$

independent of  $t$ . This is natural if  $u$  is a probability distribution, so that the constant is 1. The expected value evolution equation (8.9) has a *maximum principle*

$$\max_{x \in \mathcal{S}} f(x, t) \leq \max_{x \in \mathcal{S}} f(x, t+1) . \quad (8.14)$$

This is a natural consequence of the interpretation of  $f$  as an expectation value. The probabilities,  $u(x, t)$  need not satisfy a maximum principle either forward or backward in time.

This duality relation has is particularly transparent in matrix terms. The formula (8.8) is expressed explicitly in terms of the probabilities at time  $t$  as

$$\sum_{x \in \mathcal{S}} f(x, T)u(x, T) ,$$

which has the matrix form

$$\underline{u}(T)\underline{f}(T) .$$

Written in this order, the matrix multiplication is compatible; the other order,  $\underline{f}(T)\underline{u}(T)$ , would represent an  $n \times n$  matrix instead of a single number. In view of (8.7), we may rewrite this as

$$\underline{u}(0)P^T\underline{f}(T) .$$

Because matrix multiplication is associative, this may be rewritten

$$[\underline{u}(0)P^t] \cdot [P^{T-t}\underline{f}(T)] \quad (8.15)$$

for any  $t$ . This is the same as saying that  $\underline{u}(t)\underline{f}(t)$  is independent of  $t$ , as we already saw.

In linear algebra and functional analysis, “adjoint” or “dual” is a fancy generalization of the transpose operation of matrices. People who don’t like to think of putting the vector to the left of the matrix think of  $\underline{u}P$  as multiplication of (the transpose of)  $\underline{u}$ , on the right, by the transpose (or adjoint or dual) of  $P$ . In other words, we can do enough evolution to compute an expected value either using  $P$  its dual (or adjoint or transpose). This is the origin of the term “duality” in this context.

**Exercise 8.1** (European options as a Markov chain). Consider the case with interest rate  $r = 0$ . Then the finite difference method in Example 6.1 for a European option takes the form

$$\begin{aligned} \bar{f}_{n-1,i} &= \bar{f}_{n,i}(1 - \sigma^2 i^2 \Delta t) + \frac{1}{2} \sigma^2 i^2 \bar{f}_{n,i+1} \Delta t \\ &+ \frac{1}{2} \sigma^2 i^2 \bar{f}_{n,i-1} \Delta t, \end{aligned}$$

which is a Markov chain model called the trinomial tree method. Identify the transition probabilities.

## 8.5 Dynamic Programming

Dynamic programming is a method for valuing American style options and other financial instruments that allow the holder to make decisions that effect the ultimate payout. The idea is to define the appropriate value function,  $f(x, t)$ , that satisfies a nonlinear version of the backwards evolution equation (8.9). In the real world, dynamic programming is used to determine “optimal” trading strategies for traders trying to take or unload a big position without moving the market, to find cost efficient hedging strategies when trading costs or other market frictions are significant, and for many other purposes. Its main drawback stems from the necessity of computing the cost to go function (see below) for every state  $x \in \mathcal{S}$ . For complex models, the state space may be too large for this to be practical. That’s when things really get interesting.

I will explain the idea in a simple but somewhat abstract situation. As in the previous section, it is possible to use these ideas to treat other related problems. We have a Markov chain as before, but now the transition probabilities depend on a “control parameter”,  $\xi$ . That is

$$p(x, y, \xi) = \mathbf{Pr}(X(t+1) = y | X(t) = x, \xi) \ .$$

In the “stochastic control problem”, we are allowed to choose the control parameter at time  $t$ ,  $\xi(t)$ , knowing the value of  $X(t)$  but not any more about the future than the transition probabilities. Because the system is a Markov chain, knowledge of earlier values,  $X(t-1), \dots$ , will not help predict or control the future. Choosing  $\xi$  as a function

of  $X(t)$  and  $t$  is called “feedback control” or a “decision strategy”. The point here is that the optimal control policy is a feedback control. That is, instead of trying to choose a whole control trajectory,  $\xi(t)$  for  $t = 0, 1, \dots, T$ , we instead try to choose the feedback functions  $\xi(X(t), t)$ . We will write  $\xi(X, t)$  for such a decision strategy.

Any given strategy has an expected payout, which we write

$$\mathbf{E}_\xi [V(X(T))] \text{ .}$$

Our object is to compute the value of the financial instrument under the optimal decision strategy:

$$\max_{\xi} \mathbf{E}_\xi [V(X(T))] \text{ ,} \quad (8.16)$$

and the optimal strategy that achieves this.

The appropriate collection of values for this is the “cost to go” function

$$\begin{aligned} f(x, t) &= \max_{\xi} \mathbf{E}_\xi [V(X(T)) | X(t) = x] \\ &= \max_{\xi_t} \max_{\xi_{t+1}, \xi_{t+2}, \dots, \xi_T} \mathbf{E}_\xi [V(X(T)) | X(t+1) = y] P(x, y, \xi_t) \\ &= \max_{\xi(t)} \sum_{y \in \mathcal{S}} f(y, t+1) p(x, y, \xi(t)) \text{ .} \end{aligned} \quad (8.17)$$

As before, we have “initial data”  $f(x, T) = V(x)$ . We need to compute the values  $f(x, t)$  in terms of already computed values  $f(x, t+1)$ . For this, we suppose that the optimal decision strategy at time  $t$  is not yet known but those at later times are already computed. If we use control variable  $\xi(t)$  at time  $t$ , and the optimal control thereafter, we get payout depending on the state at time  $t+1$ :

$$\mathbf{E} [f(X(t+1), t+1) | X(t) = x, \xi(t)] = \sum_{y \in \mathcal{S}} f(y, t+1) p(x, y, \xi(t)) \text{ .}$$

Maximizing this expected payout over  $\xi(t)$  gives the optimal expected payout at time  $t$ :

$$f(x, t) = \max_{\xi(t)} \sum_{y \in \mathcal{S}} f(y, t+1) p(x, y, \xi(t)) \text{ .} \quad (8.18)$$

This is the principle of dynamic programming. We replace the “multiperiod optimization problem” (8.17) with a sequence of hopefully simpler “single period” optimization problems (8.18) for the cost to go function.

**Exercise 8.2** (American options as a controlled Markov chain). Consider the case with interest rate  $r = 0$ . Then the finite difference method in Example 6.1 for an American option takes the form

$$\bar{f}_{n-1, i} = \max \left( \bar{f}_{n, i} (1 - \sigma^2 i^2 \Delta t) + \frac{1}{2} \sigma^2 i^2 \bar{f}_{n, i+1} \Delta t + \frac{1}{2} \sigma^2 i^2 \bar{f}_{n, i-1} \Delta t, \max(K - i \Delta S, 0) \right)$$

which is a dynamic programming Markov chain model. Identify the transition probabilities and the control function.

## 8.6 Examples and Exercises

1. A stationary Markov chain has three states, called  $A$ ,  $B$ , and  $C$ . The probability of going from  $A$  to  $B$  in one step is .6. The probability of staying at  $A$  is .4. The probability of going from  $B$  to  $A$  is .3. The probability of staying at  $B$  is .2, and the probability of going to  $C$  is .5. From state  $C$ , the probability of going to  $B$  is .8 and the probability of going to  $A$  is zero. The payout for state  $A$  is 1, for state  $B$  is 4, and for state  $C$  is 9.
  - a. Compute the probabilities that the system will be in state  $A$ ,  $B$ , or  $C$  after two steps, starting from state  $A$ . Use these three numbers to compute the expected payout after two steps starting from state  $A$ .
  - b. Compute the expected payouts in one step starting from state  $A$  and from state  $B$ . These are  $f(A, 1)$  and  $f(B, 1)$  respectively.
  - c. See that the appropriate average of  $f(A, 1)$  and  $f(B, 1)$  agrees with the answer from part a.
2. Suppose a stock price is a stationary Markov chain with the following transition probabilities. In one step, the stock goes from  $S$  to  $uS$  with probability  $p$  and from  $S$  to  $dS$  with probability  $q = 1 - p$ . We generally suppose that  $u$  (the uptick) is slightly bigger than one while  $d$  (the downtick) as a bit smaller. Show that the method for computing the expected payout is exactly the binomial tree method for valuing European style options.
3. Formulate the American style option valuation problem as an optimal decision problem. Choosing the early exercise time is the same as deciding on each day whether to exercise or not. Show that the dynamic programming algorithm discussed above is the binomial tree method for American style options. The optimization problem (8.18) reduces to taking the max between the computed  $f$  and the intrinsic value.
4. This is the simplest example of the “linear quadratic gaussian” (LQG) paradigm in optimal control that has become the backbone of traditional control engineering. Here  $X(t)$  is a real number. The transitions are given by

$$X(t + 1) = aX(t) + \sigma G(t) + \xi(t) , \quad (8.19)$$

where  $G(t)$  is a standard normal random variable and the  $G(t)$  for different  $t$  values are independent. We want to minimize the quantity

$$C = \sum_{t=1}^T X(t)^2 + \mu \sum_{t=0}^{T-1} \xi(t)^2 \quad (8.20)$$

We want to find a choice of the control,  $\xi$ , that minimizes  $\mathbf{E}(C)$ . Note that the dynamics (8.19) are linear, the noise is gaussian, and the cost function (8.20) is quadratic. Define the cost to go function  $f(x, t)$  to be the cost incurred starting at

$x$  at time  $t$  ignoring the costs that are incurred at earlier times. Start by computing  $f(x, T - 1)$  explicitly by minimizing over the single variable  $\xi(T - 1)$ . Note that the optimal  $\xi(T - 1)$  is a linear function of  $X(T - 1)$ . Using this information, compute  $f(x, T - 2)$  by optimizing over  $\xi(T - 2)$ , and so on. The LQG model in control engineering justifies linear feedback control in much the same way the gaussian error model and maximum likelihood justifies least squares estimation in statistics.

## Chapter 9

# Optimal Control and Inverse Problems

The purpose of Optimal Control is to influence the behavior of a dynamical system in order to achieve a desired goal. *Optimal control* has a large variety of applications where the dynamics can be controlled optimally, such as aerospace, aeronautics, chemical plants, mechanical systems, finance and economics, but also to solve *inverse problems* where the goal is to determine input data in an equation from its solution values. An important application we will study in several settings is to determine the "data" in differential equations models using optimally controlled reconstructions of measured "solution" values.

Inverse problems are typically harder to solve numerically than forward problems since they are often ill-posed (in contrast to forward problems), where ill-posed is the opposite of well-posed and a problem is defined to be well-posed if the following three properties holds

- (1) there is a solution,
- (2) the solution is unique, and
- (3) the solution depends continuously on the data.

It is clear that a solution that does not depend continuously on its data is difficult to approximate accurately, since a tiny perturbation of the data (either as measurement error and/or as numerical approximation error) may give a large change in the solution. Therefore, the ill-posedness of inverse and optimal control problems means that they need to be somewhat modified to be solved: we call this to *regularize* the problem. Optimal control theory is suited to handle many inverse problems for differential equations, since we may formulate the objective – for instance to optimally reconstruct measured data or to find an optimal design – with the differential equation as a constraint. This chapter explains:

- the reason to regularize inverse problems in an optimal control setting,



- a method how to regularize the control problem, and
- in what sense the regularized problem approximates the original problem.

To give some intuition on optimal control and to introduce some basic concepts let us consider a hydro-power generator in a river. Suppose that we are the owners of such a generator, and that our goal is to maximise our profit by selling electricity in some local electricity market. This market will offer us buying prices at different hours, so one decision we have to make is when and how much electricity to generate. To make this decision may not be a trivial task, since besides economic considerations, we also have to meet technical constraints. For instance, the power generated is related to the amount of water in the reservoir, the turbined flow and other variables. Moreover, if we want a plan for a period longer than just a few days the water inflow to the lake may not be precisely known, making the problem stochastic.

We can state our problem in optimal control terms as the maximization of an *objective function*, the expected profit from selling electricity power during a given period, with respect to *control functions*, like the hourly turbined flow. Observe that the turbined flow is positive and smaller than a given maximum value, so it is natural to have a set of *feasible controls*, namely the set of those controls we can use in practice. In addition, our dynamical system evolves according to a given law, also called *the dynamics*, which here comes from a mass balance in the dam's lake. This law tells us how the *state variable*, the amount of water in the lake, evolves with time according to the control we give. Since the volume in the lake cannot be negative, there exist additional constraints, known as *state constraints*, that have to be fulfilled in the optimal control problem.

After introducing the formulation of an optimal control problem the next step is to find its solution. As we shall see, the optimal control is closely related with the solution of a nonlinear partial differential equation, known as the Hamilton-Jacobi-Bellman equation. To derive the Hamilton-Jacobi-Bellman equation we shall use the dynamic programming principle, which relates the solution of a given optimal control problem with solutions to simpler problems.

## 9.1 The Deterministic Optimal Control Setting

A mathematical setting for optimally controlling the solution to a deterministic ordinary differential equation

$$\begin{aligned} \dot{X}^s &= f(X^s, \alpha^s) \quad t < s < T \\ X^t &= x \end{aligned} \tag{9.1}$$

is to minimize

$$\inf_{\alpha \in \mathcal{A}} \left( \int_t^T h(X^s, \alpha^s) ds + g(X^T) \right) \tag{9.2}$$

for given *cost functions*  $h : \mathbb{R}^d \times [t, T] \rightarrow \mathbb{R}$  and  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  and a given set of *control functions*  $\mathcal{A} = \{\alpha : [t, T] \rightarrow A\}$  and *flux*  $f : \mathbb{R}^d \times A \rightarrow \mathbb{R}^d$ . Here  $A$  is a given compact subset of some  $\mathbb{R}^m$ .

### 9.1.1 Examples of Optimal Control

**Example 9.1** (Optimal control of spacecraft). To steer a spacecraft with minimal fuel consumption to an astronomical body may use the gravitational force from other bodies. The dynamics is determined by the classical Newton's laws with forces depending on the gravity on the spacecraft and its rocket forces, which is the control cf. [?].

**Example 9.2** (Inverse problem: Parameter reconstruction). The option values can be used to determine the volatility function implicitly. The objective in the optimal control formulation is then to find a volatility function that yields option prices that deviate as little as possible from the measured option prices. The dynamics is the Black-Scholes equation with the volatility function to be determined, that is the dynamics is a deterministic partial differential equation and the volatility is the control function, see Section 9.2.1.1. This is a typical inverse problem: it is called inverse because in the standard view of the Black-Scholes equation relating the option values and the volatility, the option price is the unknown and the volatility is the data; while here the formulation is reversed with option prices as data and volatility as unknown in the same Black-Scholes equation.

**Example 9.3** (Inverse problem: Weather prediction). The incompressible Navier-Stokes equations are used to forecast weather. The standard mathematical setting of this equation is an initial value problem with unknown velocity and pressure to be determined from the initial data: in weather prediction one can use measured velocity and pressure not only at a single initial instance but data given over a whole time history. An optimal control formulation of the weather prediction is to find the first initial data (the control) matching the time history of measured velocity and pressure with the Navier-Stokes dynamics as constraint. Such an optimal control setting improves the accuracy and makes longer forecast possible as compared to the classical initial value problem, see [Pir84], [?]. This is an inverse problem since the velocity and pressure are used to determine the "initial data".

**Example 9.4** (Merton's stochastic portfolio problem). A basic problem in finance is to choose how much to invest in stocks and in bonds to maximize a final utility function. The dynamics of the portfolio value is then stochastic and the objective is to maximize an expected value of a certain (utility) function of the portfolio value, see section 9.3.1.

**Example 9.5** (Euler-Lagrange equation). The shape of a soap bubble between a wire frame can be determined as the surface that minimizes the bubble area. For a surface in  $\mathbb{R}^3$  described by  $\{(x, u(x)) : x \in \Omega \subset \mathbb{R}^2\}$  the area is given by

$$\int_{\Omega} \sqrt{1 + |\nabla u|^2} dx.$$

Here the whole surface is the control function, and given a wire  $\{(x, g(x)) : x \in \partial\Omega\}$ , the minimal surface solves the Euler-Lagrange equation,

$$\begin{aligned} \operatorname{div} \left( \frac{\nabla u}{\sqrt{1 + |\nabla u|^2}} \right) &= 0, \quad \text{in } \Omega, \\ u &= g, \quad \text{on } \partial\Omega. \end{aligned}$$

**Example 9.6** (Inverse problem: Optimal design). An example of optimal design is to construct an electrical conductor to minimize the power loss by placing a given amount of conductor in a given domain, see Section 9.2.1.2. This is an inverse problem since the conductivity is determined from the electric potential in an equation where the standard setting is to determine the electric potential from the given conductivity.

### 9.1.2 Approximation of Optimal Control

Optimal control problems can be solved by the Lagrange principle or dynamic programming. The *dynamic programming* approach uses the value function, defined by

$$u(x, t) := \inf_{\alpha \in \mathcal{A}} \left( \int_t^T h(X^s, \alpha^s) ds + g(X^T) \right), \quad (9.3)$$

for the ordinary differential equation (9.1) with  $X_t \in \mathbb{R}^d$ , and leads to solution of a non linear *Hamilton-Jacobi-Bellman* partial differential equation

$$\partial_t u(x, t) + \underbrace{\min_{\alpha \in \mathcal{A}} (f(x, \alpha) \cdot \partial_x u(x, t) + h(x, \alpha))}_{H(\partial_x u(x, t), x)} = 0, \quad t < T, \quad (9.4)$$

$$u(\cdot, T) = g,$$

in  $(x, t) \in \mathbb{R}^d \times \mathbb{R}_+$ . The *Lagrange principle* (which seeks a minimum of the cost with the dynamics as a constraint) leads to the solution of a Hamiltonian system of ordinary differential equations, which are the characteristics of the Hamilton-Jacobi-Bellman equation

$$\begin{aligned} X'^t &= f(X^t, \alpha^t), \quad X_0 \text{ given,} \\ -\lambda'_i &= \partial_{x_i} f(X^t, \alpha^t) \cdot \lambda^t + \partial_{x_i} h(X^t, \alpha^t), \quad \lambda^T = g'(X^T), \\ \alpha^t &\in \operatorname{argmin}_{a \in \mathcal{A}} \left( \lambda^t \cdot f(X^t, a) + h(X^t, a) \right), \end{aligned} \quad (9.5)$$

based on the Pontryagin Principle. The next sections explain these two methods.

The non linear Hamilton-Jacobi partial differential approach has the theoretical advantage of well established theory and that a global minimum is found; its fundamental drawback is that it cannot be used computationally in high dimension  $d \gg 1$ , since the computational work increases exponentially with the dimension  $d$ . The Lagrange principle has the computational advantage that high dimensional problems,  $d \gg 1$ , can often be solved and its drawback is that in practice only local minima can be found computationally, often with some additional error introduced by a regularization method. Another drawback with the Lagrange principle is that it (so far) has no efficient implementation in the natural stochastic setting with adapted Markov controls, while the Hamilton-Jacobi PDE approach directly extends to such stochastic controls, see Section 9.3; as a consequence computations of stochastic controls is basically limited to low dimensional problems.

### 9.1.3 Motivation of the Lagrange formulation

Let us first review the Lagrange multiplier method to minimize a function subject to a constraint  $\min_{x \in A, y=g(x)} F(x, y)$ . Assume  $F : \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}$  is a differentiable function. The goal is to find the minimum  $\min_{x \in A} F(x, g(x))$  for a given differentiable function  $g : \mathbb{R}^d \rightarrow \mathbb{R}^n$  and a compact set  $A \subset \mathbb{R}^d$ . This problem leads to the usual necessary condition for an interior minimum

$$\frac{d}{dx} F(x, g(x)) = \partial_x F(x, g(x)) + \partial_y F(x, g(x)) \partial_x g(x) = 0. \quad (9.6)$$

An alternative method to find the solution is to introduce the Lagrangian function  $\mathcal{L}(\lambda, y, x) := F(x, y) + \lambda \cdot (y - g(x))$  with the Lagrange multiplier  $\lambda \in \mathbb{R}^n$  and choose  $\lambda$  appropriately to write the necessary condition for an interior minimum

$$\begin{aligned} 0 &= \partial_\lambda \mathcal{L}(\lambda, y, x) = y - g(x), \\ 0 &= \partial_y \mathcal{L}(\lambda, y, x) = \partial_y F(x, y) + \lambda, \\ 0 &= \partial_x \mathcal{L}(\lambda, y, x) = \partial_x F(x, y) - \lambda \cdot \partial_x g(x). \end{aligned}$$

Note that the first equation is precisely the constraint. The second equation determines the multiplier to be  $\lambda = -\partial_y F(x, y)$ . The third equation yields for this multiplier  $\partial_x \mathcal{L}(-\partial_y F(x, y), y, x) = \frac{d}{dx} F(x, g(x))$ , that is the multiplier is chosen precisely so that the partial derivative with respect to  $x$  of the Lagrangian is the total derivative of the objective function  $F(x, g(x))$  to be minimized. This Lagrange principle is often practical to use when the constraint is given implicitly, e.g. as  $g(x, y) = 0$  with a differentiable  $g : \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}^n$ ; then the condition  $\det \partial_y g(x, y) \neq 0$  in the implicit function theorem implies that the function  $y(x)$  is well defined and satisfies  $g(x, y(x)) = 0$  and  $\partial_x y = -\partial_y g(x, y)^{-1} \partial_x g(x, y)$ , so that the Lagrange multiplier method works.

The Lagrange principle for the optimal control problem (9.1) -(9.2), to minimize the cost with the dynamics as a constraint, leads to the Lagrangian

$$\mathcal{L}(\lambda, X, \alpha) := g(X^T) + \int_0^T h(X^s, \alpha^s) ds + \int_0^T \lambda^s \cdot (f(X^s, \alpha^s) - \dot{X}) ds \quad (9.7)$$

with a Lagrange multiplier function  $\lambda : [0, T] \rightarrow \mathbb{R}^d$ . Differentiability of the Lagrangian leads to the necessary conditions for a constrained interior minimum

$$\begin{aligned} \partial_\lambda \mathcal{L}(X, \lambda, \alpha) &= 0, \\ \partial_X \mathcal{L}(X, \lambda, \alpha) &= 0, \\ \partial_\alpha \mathcal{L}(X, \lambda, \alpha) &= 0. \end{aligned} \quad (9.8)$$

Our next step is to verify that the two first equations above are the same as the two first in (9.5) and that the last equation is implied by the stronger Pontryagin principle in the last equation in (9.5). We will later use the Hamilton-Jacobi equation in the dynamic programming approach to verify the Pontryagin principle.

*The first equation.* Choose a real valued continuous function  $v : [0, T] \rightarrow \mathbb{R}^d$  and define the function  $L : \mathbb{R} \rightarrow \mathbb{R}$  by  $L(\epsilon) := \mathcal{L}(X, \lambda + \epsilon v, \alpha)$ . Then the first of the three equations means precisely that  $L'(0) = \frac{d}{d\epsilon} \mathcal{L}(X, \lambda + \epsilon v, \alpha)|_{\epsilon=0} = 0$ , which implies that

$$0 = \int_0^T v^s \cdot (f(X^s, \alpha^s) - \dot{X}^s) ds$$

for any continuous function  $v$ . If we assume that  $f(X^s, \alpha^s) - \dot{X}^s$  is continuous we obtain  $f(X^s, \alpha^s) - \dot{X}^s = 0$ : since if  $\beta(s) := f(X^s, \alpha^s) - \dot{X}^s \neq 0$  for some  $s$  there is an interval where  $\beta$  is either positive or negative; by choosing  $v$  to be zero outside this interval we conclude that  $\beta$  is zero everywhere and we have derived the first equation in (9.5).

*The second equation.* The next equation  $\frac{d}{d\epsilon} \mathcal{L}(X + \epsilon v, \lambda, \alpha)|_{\epsilon=0} = 0$  needs  $v^0 = 0$  by the initial condition on  $X^0$  and leads by integration by parts to

$$\begin{aligned} 0 &= \int_0^T \lambda^s \cdot (\partial_{X_i} f(X^s, \alpha^s) v_i^s - \dot{v}^s) + \partial_{X_i} h(X^s, \alpha^s) v_i^s ds + \partial_{X_i} g(X^T) v_i^T \\ &= \int_0^T \lambda^s \cdot \partial_{X_i} f(X^s, \alpha^s) v_i^s + \dot{\lambda} \cdot v^s + \partial_{X_i} h(X^s, \alpha^s) v_i^s ds \\ &\quad + \lambda^0 \cdot \underbrace{v^0}_{=0} - (\lambda^T - \partial_X g(X^T)) \cdot v^T \\ &= \int_0^T (\partial_X f^*(X^s, \alpha^s) \lambda^s + \dot{\lambda}^s + \partial_X h(X^s, \alpha^s)) \cdot v^s ds \\ &\quad - (\lambda^T - \partial_X g(X^T)) \cdot v^T, \end{aligned}$$

using the summation convention  $a_i b_i := \sum_i a_i b_i$ . Choose now the function  $v$  to be zero outside an interior interval where possibly  $\partial_X f^*(X^s, \alpha^s) \lambda^s + \dot{\lambda}^s + \partial_X h(X^s, \alpha^s)$  is non zero, so that in particular  $v^T = 0$ . We see then that in fact  $\partial_X f^*(X^s, \alpha^s) \lambda^s + \dot{\lambda}^s + \partial_X h(X^s, \alpha^s)$  must be zero (as for the first equation) and we obtain the second equation in (9.5). Since the integral in the right hand side vanishes, varying  $v^T$  shows that the final condition for the Lagrange multiplier  $\lambda^T - \partial_X g(X^T) = 0$  also holds.

*The third equation.* The third equation in (9.8) implies as above that for any function  $v(t)$  compactly supported in  $A$

$$0 = \int_0^T \lambda^s \cdot \partial_\alpha f(X^s, \alpha^s) v + \partial_\alpha h(X^s, \alpha^s) v ds$$

which yields

$$\lambda^s \cdot \partial_\alpha f(X^s, \alpha^s) + \partial_\alpha h(X^s, \alpha^s) = 0 \tag{9.9}$$

in the interior  $\alpha \in A - \partial A$  minimum point  $(X, \lambda, \alpha)$ . The last equation in (9.5) is a stronger condition: it says that  $\alpha$  is a minimizer of  $\lambda^s \cdot f(X^s, a) + h(X^s, a) = 0$  with respect to  $a \in A$ , which clearly implies (9.9) for interior points  $\alpha \in A - \partial A$ . To derive the Pontryagin principle we will use dynamic programming and the Hamilton-Jacobi-Bellman equation which is the subject of the next section.

### 9.1.4 Dynamic Programming and the Hamilton-Jacobi-Bellman Equation

The dynamic programming view to solve optimal control problems is based on the idea to track the optimal solution backwards: at the final time the value function is given  $u(x, T) = g(x)$  and then, recursively for small time step backwards, find the optimal control to go from each point  $(x, t)$  on the time level  $t$  to the time level  $t + \Delta t$  with the value function  $u(\cdot, t + \Delta t)$ , see Figure 9.1. Assume for simplicity first that  $h \equiv 0$  then any path  $X : [t, t + \Delta t] \rightarrow \mathbb{R}^d$  starting in  $X^t = x$  will satisfy

$$u(x, t) = \inf_{\alpha: [t, t+\Delta t] \rightarrow A} u(X^{t+\Delta t}, t + \Delta t),$$

so that if  $u$  is differentiable

$$du(X^t, t) = (\partial_t u(X^t, t) + \partial_x u(X^t, t) \cdot f(X^t, \alpha_t)) dt \geq 0, \quad (9.10)$$

since a path from  $(x, t)$  with value  $u(x, t)$  can lead only to values  $u(X^{t+\Delta t}, t + \Delta t)$  which are not smaller than  $u(x, t)$ . If also the infimum is attained, then an optimal path  $X_*^t$  exists, with control  $\alpha_*^t$ , and satisfies

$$du(X_*^t, t) = (\partial_t u(X_*^t, t) + \partial_x u(X_*^t, t) \cdot f(X_*^t, \alpha_*^t)) dt = 0. \quad (9.11)$$

The combination of (9.10) and (9.11) implies that

$$\begin{aligned} \partial_t u(x, t) + \min_{\alpha \in A} (\partial_x u(x, t) \cdot f(x, \alpha)) &= 0 \quad t < T \\ u(\cdot, T) &= g, \end{aligned}$$

which is the Hamilton-Jacobi-Bellman equation in the special case  $h \equiv 0$ .

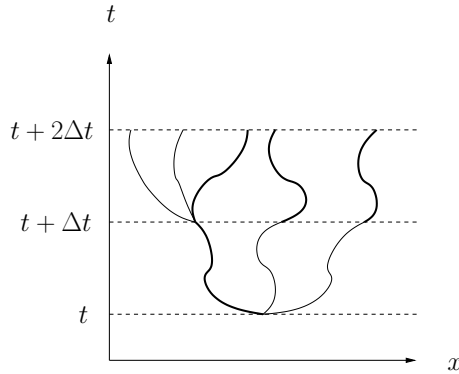


Figure 9.1: Illustration of dynamics programming.

The case with  $h$  non zero follows similarly by noting that now

$$0 = \inf_{\alpha: [t, t+\Delta t] \rightarrow A} \left( \int_t^{t+\Delta t} h(X^s, \alpha^s) ds + u(X^{t+\Delta t}, t + \Delta t) - u(x, t) \right), \quad (9.12)$$

which for differentiable  $u$  implies the Hamilton-Jacobi-Bellman equation (9.4)

$$\begin{aligned} 0 &= \inf_{\alpha \in A} (h(x, \alpha) + \partial_t u(x, t) + \partial_x u(x, t) \cdot f(x, \alpha)) \\ &= \partial_t u(x, t) + \underbrace{\min_{\alpha \in A} (\partial_x u(x, t) \cdot f(x, \alpha) + h(x, \alpha))}_{=: H(\partial_x u(x, t), x)} \quad t < T, \\ g &= u(\cdot, T). \end{aligned}$$

Note that this derivation did not assume that an optimal path is attained, but that  $u$  is differentiable which in general is not true. There is fortunately a complete theory for non differentiable solutions to Hamilton-Jacobi equations, with its basics presented in Section 9.1.6. First we shall relate the Lagrange multiplier method with the Pontryagin principle to the Hamilton-Jacobi-Bellman equation using characteristics.

### 9.1.5 Characteristics and the Pontryagin Principle

The following theorem shows that the characteristics of the Hamilton-Jacobi equation is a Hamiltonian system.

**Theorem 9.7.** *Assume  $u \in C^2$ ,  $H \in C^1$  and*

$$\dot{X}^t = \partial_\lambda H(\lambda^t, X^t)$$

with  $\lambda^t := \partial_x u(X^t, t)$ . Then the characteristics  $(X^t, \lambda^t)$  satisfy the Hamiltonian system

$$\begin{aligned} \dot{X}^t &= \partial_\lambda H(\lambda^t, X^t) \\ \dot{\lambda}^t &= -\partial_X H(\lambda^t, X^t). \end{aligned} \tag{9.13}$$

*Proof.* The goal is to verify that the construction of  $X^t$  implies that  $\lambda$  has the dynamics (9.13). The definition  $\dot{X}^t = \partial_\lambda H(\lambda^t, X^t)$  implies by  $x$ -differentiation of the Hamilton-Jacobi equation along the path  $(X^t, t)$

$$\begin{aligned} 0 &= \partial_{x_k} \partial_t u(X^t, t) + \sum_j \partial_{\lambda_j} H(\partial_x u(X^t, t), X^t) \underbrace{\partial_{x_k} \partial_{x_j} u(X^t, t)}_{= \partial_{x_j} \partial_{x_k} u} + \partial_{x_k} H(\partial_x u(X^t, t), X^t) \\ &= \frac{d}{dt} \partial_{x_k} u(X^t, t) + \partial_{x_k} H(\partial_x u(X^t, t), X^t) \end{aligned}$$

which by the definition  $\lambda^t := \partial_x u(X^t, t)$  is precisely  $\dot{\lambda}^t + \partial_x H(\lambda^t, X^t) = 0$ .  $\square$

The next step is to relate the characteristics  $X^t, \lambda^t$  to the solution of the Lagrange principle (9.5). But note first that the Hamiltonian  $H$  in general is not differentiable, even if  $f$  and  $h$  are very regular: for instance  $\dot{X} = f(X^t)$  and  $h(x, \alpha) = x\alpha$  implies for  $A = [-1, 1]$  that the Hamiltonian becomes  $H(\lambda, x) = \lambda f(x) - |x|$  which is only Lipschitz continuous, that is  $|H(\lambda, x) - H(\lambda, y)| \leq K|x - y|$  with the Lipschitz constant  $K = 1 + \|\lambda \cdot \partial_x f(\cdot)\|_\infty$  in this case. In fact if  $f$  and  $h$  are bounded differentiable functions the Hamiltonian will always be Lipschitz continuous satisfying  $|H(\lambda, x) - H(\nu, y)| \leq K(|\lambda - \nu| + |x - y|)$  for some constant  $K$ , see Exercise ??.

**Theorem 9.8.** *Assume that  $f$ ,  $h$  are  $x$ -differentiable in  $(x, \alpha_*)$  and a control  $\alpha_*$  is optimal for a point  $(x, \lambda)$ , i.e.*

$$\lambda \cdot f(x, \alpha_*) + h(x, \alpha_*) = H(\lambda, x),$$

*and suppose also that  $H$  is differentiable in the point or that  $\alpha_*$  is unique. Then*

$$\begin{aligned} f(x, \alpha_*) &= \partial_\lambda H(\lambda, x), \\ \lambda \cdot \partial_{x_i} f(x, \alpha_*) + \partial_{x_i} h(x, \alpha_*) &= \partial_{x_i} H(\lambda, x). \end{aligned} \tag{9.14}$$

*Proof.* We have for any  $w, v \in \mathbb{R}^d$

$$\begin{aligned} H(\lambda + w, x + v) - H(\lambda, x) &\leq (\lambda + w) \cdot f(x + v, \alpha_*) + h(x + v, \alpha_*) \\ &\quad - \lambda \cdot f(x, \alpha_*) - h(x, \alpha_*) \\ &= w \cdot f(x, \alpha_*) + \sum_{i=1}^d (\lambda \cdot \partial_{x_i} f + \partial_{x_i} h) v_i + o(|v| + |w|) \end{aligned}$$

which implies (9.14) by choosing  $w$  and  $v$  in all directions.  $\square$

This Theorem shows that the Hamiltonian system (9.13) is the same as the system (9.5), given by the Lagrange principle using the optimal control  $\alpha_*$  with the Pontryagin principle

$$\lambda \cdot f(x, \alpha_*) + h(x, \alpha_*) = \inf_{\alpha \in A} (\lambda \cdot f(x, \alpha) + h(x, \alpha)) =: H(\lambda, x).$$

If  $\alpha_*$  is not unique (i.e not a single point) the proof shows that (9.14) still holds for the optimal controls, so that  $\partial_\lambda H$  and  $\partial_x H$  become set valued. We conclude that non unique local controls  $\alpha_*$  is the phenomenon that makes the Hamiltonian non differentiable in certain points. In particular a differentiable Hamiltonian gives unique optimal control fluxes  $\partial_\lambda H$  and  $\partial_x H$ , even if  $\alpha_*$  is not a single point. If the Hamiltonian can be explicitly formulated, it is therefore often practical to use the Hamiltonian system formulation with the variables  $X$  and  $\lambda$ , avoiding the control variable.

Clearly, the Hamiltonian needs to be differentiable for the Hamiltonian system to make sense; in fact its flux  $(\partial_\lambda H, -\partial_x H)$  must be Lipschitz continuous to give well posedness. On the other hand we shall see that the Hamilton-Jacobi-Bellman formulation, based on dynamic programming, leads to non differentiable value functions  $u$ , so that classical solutions lack well posedness. The mathematical setting for optimal control therefore seemed somewhat troublesome both on the Hamilton-Jacobi PDE level and on the Hamilton ODE level. In the 1980's the situation changed: Crandall-Lions-Evans [CEL84] formulated a complete well posedness theory for generalized so called viscosity solutions to Hamilton-Jacobi partial differential equations, allowing Lipschitz continuous Hamiltonians. The theory of viscosity solutions for Hamilton-Jacobi-Bellman partial differential equations provides good theoretical foundation also for non smooth controls.



In particular this mathematical theory removes one of Pontryagin’s two reasons<sup>1</sup>, but not the other, to favor the ODE approach (9.5) and (9.13): the mathematical theory of viscosity solutions handles elegantly the inherent non smoothness in control problems; analogous theoretical convergence results for an ODE approach was developed later based on the so called minmax solutions, see [Sub95]; we will use an alternative ODE method to solve optimal control problems numerically based on regularized Hamiltonians, where we approximate the Hamiltonian with a two times differentiable Hamiltonian, see Section 9.2.

Before we formulate the generalized solutions, we show that classical solutions only exist for short time in general.

**Example 9.9.** The Hamilton-Jacobi equation

$$\partial_t u - \frac{1}{2}(\partial_x u)^2 = 0$$

has the characteristics

$$\begin{aligned}\dot{X}^t &= -\lambda^t \\ \dot{\lambda}^t &= 0,\end{aligned}$$

which implies  $\dot{X}^t = \text{constant}$ . If the initial data  $u(\cdot, T)$  is a concave function (e.g. a smooth version of  $-|x|$ ) characteristics  $X$  will collide, see Figure 9.2. We can understand this precisely by studying blow-up of the derivative  $w$  of  $\partial_x u =: v$ ; since  $v$  satisfies

$$\partial_t v - \underbrace{\frac{1}{2}\partial_x(v^2)}_{v\partial_x v} = 0$$

we have by  $x$ -differentiation

$$\underbrace{\partial_t w - v\partial_x w - w^2}_{\frac{d}{dt}w(X^t, t)} = 0,$$

which reduces to the ordinary differential equation for  $z^t := w(X^t, t)$

$$\frac{d}{dt}z(t) = z^2(t).$$

Its separation of variables solution  $dz/z^2 = dt$  yields  $-1/z^t = t + C$ . The constant becomes  $C = -T - 1/z^T$ , so that  $z^t = 1/(t - T - 1/z^T)$  blows up to infinity at time  $T - t = 1/z^T$ . For instance if  $z^T = -10$ , the time to blow-up time is  $1/10$ .

---

<sup>1</sup> citation from chapter one in [PBG64] “This equation of Bellman’s yields an approach to the solution of the optimal control problem which is closely connected with, but different from, the approach described in this book (see Chapter 9). It is worth mentioning that the assumption regarding the continuous differentiability of the functional (9.8) [(9.3) here] is not fulfilled in even the simplest cases, so that Bellman’s arguments yield a good heuristic method rather than a mathematical solution of the problem. The maximum principle, apart from its sound mathematical basis, also has the advantage that it leads to a system of ordinary differential equations, whereas Bellman’s approach requires the solution of a partial differential equation.”

### 9.1.6 Generalized Viscosity Solutions of Hamilton-Jacobi-Bellman Equations

Example 9.9 shows that Hamilton-Jacobi equations do in general not have global classical solutions – after finite time the derivative can become infinitely large even with smooth initial data and a smooth Hamiltonian. Therefore a more general solution concept is needed. We shall describe the so called viscosity solutions introduced by Crandall and Lions in [?], which can be characterised by the limit of viscous approximations  $u^\epsilon$  satisfying for  $\epsilon > 0$

$$\begin{aligned} \partial_t u^\epsilon(x, t) + H(\partial_x u^\epsilon(x, t), x) + \epsilon \partial_{xx} u^\epsilon(x, t) &= 0 \quad t < T \\ u^\epsilon(\cdot, T) &= g. \end{aligned}$$

The function  $u^\epsilon$  is also a value function, now for the stochastic optimal control problem

$$dX^t = f(X^t, \alpha^t)dt + \sqrt{2\epsilon} dW^t \quad t > 0$$

with the objective to minimize

$$\min_{\alpha} \mathbb{E} \left[ g(X^T) + \int_0^T h(X^t, \alpha^t) dt \mid X^0 \text{ given} \right],$$

over adapted controls  $\alpha : [0, T] \rightarrow A$ , where  $W : [0, \infty) \rightarrow \mathbb{R}^d$  is the  $d$ -dimensional Wiener process with independent components. Here adapted controls means that  $\alpha_t$  does not use values of  $W^s$  for  $s > t$ . Section 9.3 shows that the value function for this optimal control problem solves the second order Hamilton-Jacobi equation, that is

$$u^\epsilon(x, t) = \min_{\alpha} \mathbb{E} \left[ g(X^T) + \int_0^T h(X^t, \alpha^t) dt \mid X^t = x \right].$$

**Theorem 9.10** (Crandall-Lions). *Assume  $f$ ,  $h$  and  $g$  are Lipschitz continuous and bounded, then the limit  $\lim_{\epsilon \rightarrow 0+} u^\epsilon$  exists. This limit is called the viscosity solution of the Hamilton-Jacobi equation*

$$\begin{aligned} \partial_t u(x, t) + H(\partial_x u(x, t), x) &= 0 \quad t < T \\ u(\cdot, T) &= g. \end{aligned} \tag{9.15}$$

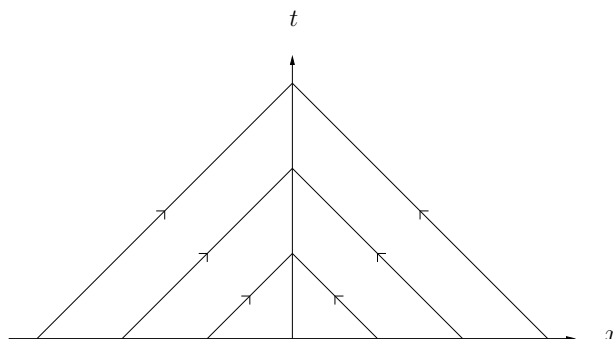


Figure 9.2: Characteristic curves colliding.

There are several equivalent ways to describe the viscosity solution directly without using viscous or stochastic approximations. We shall use the one based on sub and super differentials presented first in [CEL84]. To simplify the notation introduce first the space-time coordinate  $y = (x, t)$ , the space-time gradient  $p = (p_x, p_t) \in \mathbb{R}^{d+1}$  (related to  $(\partial_x u(y), \partial_t u(y))$ ) and write the Hamilton-Jacobi operator  $F(p, y) := p_t + H(p_x, x)$ . For a bounded uniformly continuous function  $v : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}$  define for each space-time point  $y$  its sub differential set

$$D^-v(y) = \{p \in \mathbb{R}^{d+1} : \liminf_{z \rightarrow 0} |z|^{-1}(v(y+z) - v(y) - p \cdot z) \geq 0\}$$

and its super differential set

$$D^+v(y) = \{p \in \mathbb{R}^{d+1} : \limsup_{z \rightarrow 0} |z|^{-1}(v(y+z) - v(y) - p \cdot z) \leq 0\}.$$

These two sets always exist (one may be empty), see Example 9.11; they degenerate to a single point, the space-time gradient of  $v$ , precisely if  $v$  is differentiable, that is when

$$D^-v(y) = D^+v(y) = \{p\} \iff v(y+z) - v(y) - p \cdot z = o(z).$$

**Example 9.11.** Let  $u(x) = -|x|$ , then

$$D^+u(x) = D^-u(x) = \{-\text{sgn}(x)\} \quad x \neq 0$$

$$D^-u(0) = \emptyset \quad x = 0$$

$$D^+u(0) = [-1, 1] \quad x = 0$$

see Figure 9.3.

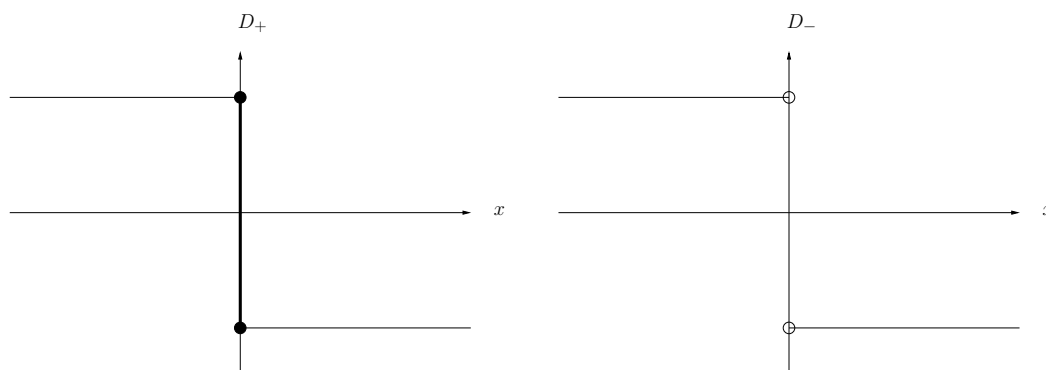


Figure 9.3: Illustration of the super and subdifferential sets for  $-|x|$ .

**Definition 9.12** (Viscosity solution). A bounded uniformly continuous function  $u$  is a viscosity solution to (9.15) if  $u(\cdot, T) = g$  and for each point  $y = (x, t)$

$$F(p, y) \geq 0 \quad \text{for all } p \in D^+u(y)$$

and

$$F(p, y) \leq 0 \quad \text{for all } p \in D^-u(y).$$

**Theorem 9.13.** *The first variation of the value function is in the superdifferential.*

*Proof.* Consider an optimal path  $X_*$ , starting in  $\bar{y} = (\bar{x}, \bar{t})$ , with control  $\alpha_*$ . We define the first variation,  $(\lambda^{\bar{t}}, \nu^{\bar{t}}) \in \mathbb{R}^d \times \mathbb{R}$ , of the value function along *this path*, with respect to perturbations in the initial point  $\bar{y}$ : let  $X_y$  be a path starting from a point  $y = (x, t)$ , close to  $\bar{y}$ , using the control  $\alpha_*$ , the differentiability of the flux  $f$  and the cost  $h$  implies that the first variation satisfies

$$\lambda_i^{\bar{t}} = \lim_{z \rightarrow 0} z^{-1} \left( \int_{\bar{t}}^T h(X_{\bar{x}+ze_i}^t, \alpha_*^t) - h(X_{\bar{x}}^t, \alpha_*^t) dt + g(X_{\bar{x}+ze_i}^T) - g(X_{\bar{x}}^T) \right) \quad (9.16)$$

and

$$\begin{aligned} -\dot{\lambda}^t &= \partial_X f(X_*^t, \alpha_*^t) \lambda^t + \partial_X h(X_*^t, \alpha_*^t) \quad \bar{t} < t < T, \\ \lambda^T &= g'(X_*^T), \end{aligned}$$

where  $e_i$  is the  $i$ th unit basis vector in  $\mathbb{R}^d$ . The definition of the value function shows that

$$-h(X_*^t, \alpha_*^t) = \frac{du}{dt}(X_*^t, t) = \lambda^t \cdot f(X_*^t, \alpha_*^t) + \nu^t$$

so that

$$\nu^t = -\lambda^t \cdot f(X_*^t, \alpha_*^t) - h(X_*^t, \alpha_*^t).$$

Since the value function is the minimum possible cost, we have by (9.16)

$$\begin{aligned} & \limsup_{s \rightarrow 0^+} s^{-1} \left( u(\bar{y} + s(y - \bar{y})) - u(\bar{y}) \right) \\ & \leq \limsup_{s \rightarrow 0^+} s^{-1} \left( \int_{\bar{t}}^T h(X_{\bar{y}+s(y-\bar{y})}^t, \alpha_*^t) dt + g(X_{\bar{y}+s(y-\bar{y})}^T) \right. \\ & \quad \left. - \int_{\bar{t}}^T h(X_{\bar{y}}^t, \alpha_*^t) dt + g(X_{\bar{y}}^T) \right) \\ & = \left( \lambda^{\bar{t}}, -(\lambda^{\bar{t}} \cdot f(X_*^{\bar{t}}, \alpha_*^{\bar{t}}) + h(X_*^{\bar{t}}, \alpha_*^{\bar{t}})) \right) \cdot (y - \bar{y}), \end{aligned}$$

which means precisely that the first variation is in the superdifferential.  $\square$

**Theorem 9.14.** *The value function is semi-concave, that is for any point  $(x, t)$  either the value function is differentiable or the sub differential is empty (i.e.  $D^-u(x, t) = \emptyset$  and  $D^+u(x, t)$  is non empty).*

*Proof.* Assume that the subdifferential  $D^-u(y)$  has at least two elements  $p_-$  and  $p_+$  (we will show that this leads to a contradiction). Then  $u$  is larger or equal to the wedge like function

$$u(y) \geq u(\bar{y}) + \max(p_- \cdot (y - \bar{y}), p_+ \cdot (y - \bar{y})), \quad (9.17)$$

see Figure 9.4. The definition of the value function shows that the right derivative satisfies

$$\limsup_{s \rightarrow 0^+} s^{-1} \left( u(\bar{y} + s(y - \bar{y})) - u(\bar{y}) \right) \leq (\lambda, \nu) \cdot (y - \bar{y}) \quad (9.18)$$

where  $(\lambda, \nu)$  is the first variation (in  $x$  and  $t$ ) of  $u$  around the optimal path starting in  $\bar{y}$ . The wedge bound (9.17) implies

$$\limsup_{s \rightarrow 0^+} s^{-1} \left( u(\bar{y} + s(y - \bar{y})) - u(\bar{y}) \right) \geq \max(p_- \cdot (y - \bar{y}), p_+ \cdot (y - \bar{y})),$$

but the value function cannot be both below a  $(\lambda, \nu)$ -half plane (9.18) and above such wedge function, see Figure 9.5. Therefore the subdifferential can contain at most one point: either the subdifferential is empty or there is precisely one point  $p$  in the subdifferential and in this case we see that the the first variation coincides with this point  $(\lambda, \nu) = p$ , that is the value function is differentiable  $\square$

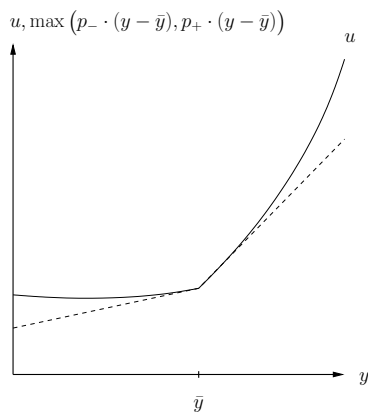


Figure 9.4: Characteristic curves colliding.

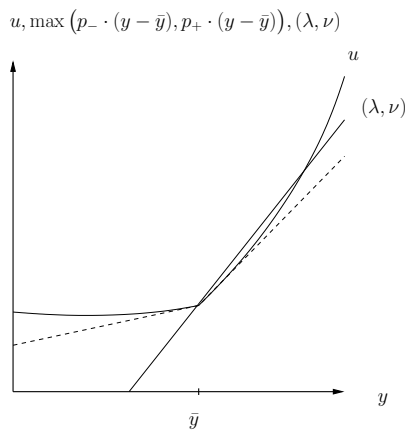


Figure 9.5: Characteristic curves colliding.

**Theorem 9.15.** *The value function is a viscosity solution.*

*Proof.* We have seen in Section 9.1.4 that for the points where the value function is differentiable it satisfies the Hamilton-Jacobi-Bellman equation. Theorem 9.14 shows that the value function  $u$  is semi-concave. Therefore, by Definition 9.12, it is enough to verify that  $p \in D^+u(x, t)$  implies  $p_t + H(p_x, x) \geq 0$ . Assume for simplicity that  $h \equiv 0$ .

There is a  $p \in D^+u(x, t)$ , which is the first variation of  $u$  along an optimal path  $(X^*, \alpha^*)$ , such that

$$\begin{aligned} p_t + H(p_x, x) &= p \cdot (f(x, \alpha), 1) \\ &\geq \limsup_{\Delta t \rightarrow 0^+} \frac{u(X^{t+\Delta t}, t + \Delta t) - u(X^t, t)}{\Delta t} = 0, \end{aligned}$$

using the definition of the superdifferential and dynamic programming. This means that any optimal control yields a super differential point  $p$  satisfying  $p_t + H(p_x, x) \geq 0$ . To finish the proof we note that any point in the super differential set can for some  $s \in [0, 1]$  be written as a convex combination  $sp^1 + (1-s)p^2$  of two points  $p^1$  and  $p^2$  in the super differential that correspond to (different) optimal controls. Since  $H$  is concave in  $p$  (see Exercise 9.19) there holds

$$\begin{aligned} &sp_t^1 + (1-s)p_t^2 + H(sp_x^1 + (1-s)p_x^2, x) \\ &\geq s(p_t^1 + H(p_x^1, x)) + (1-s)(p_t^2 + H(p_x^2, x)) \\ &\geq 0 \end{aligned}$$

which shows that  $u$  is a viscosity solution. The general case with non zero  $h$  is similar as in (9.12).  $\square$

**Theorem 9.16.** *Bounded uniformly continuous viscosity solutions are unique.*

The standard uniqueness proof uses a special somewhat complex doubling of variables technique, see [Eva98] inspired by Kruzkov. The maximum norm stability of semi-concave viscosity solutions in Section 9.1.7 also implies uniqueness.

**Example 9.17.** Consider the function  $u(x, t) = -|x|$ . We have from Example 9.11

$$D^+u(x, t) = \begin{cases} (-\operatorname{sgn}(x), 0) & x \neq 0 \\ ([-1, 1], 0) & x = 0 \end{cases}$$

and

$$D^-u(x, t) = \begin{cases} (-\operatorname{sgn}(x), 0) & x \neq 0 \\ \emptyset & x = 0. \end{cases}$$

Consequently for  $H(\lambda, x) := (1 - |\lambda|^2)/2$  we obtain

$$\begin{aligned} p_t + H(p_x, x) &\geq 0 & q \in D^+u(x, t) \\ p_t + H(p_x, x) &= 0 & q \in D^-u(x, t) \end{aligned}$$

so that  $-|x|$  is a viscosity solution to  $\partial_t u + H(\partial_x u, x) = 0$ . Similarly the function  $u(x, t) = |x|$  satisfies

$$D^-u(x, t) = \begin{cases} (\operatorname{sgn}(x), 0) & x \neq 0 \\ ([-1, 1], 0) & x = 0 \end{cases}$$

and therefore

$$p_t + H(p_x, 0) > 0 \quad \text{for } q \in (-1, 1) \subset D^-u(0, t)$$

so that  $|x|$  is not a viscosity solution to  $\partial_t u + H(\partial_x u, x) = 0$ .

### 9.1.6.1 The Pontryagin Principle for Generalized Solutions

Assume that  $X_*$  and  $\alpha_*$  is an optimal control solution. Let

$$\begin{aligned} -\dot{\lambda}_*^t &= \partial_X f(X_*^t, \alpha_*^t) \lambda_*^t + \partial_X h(X_*^t, \alpha_*^t) \quad t < T, \\ \lambda_*^T &= g'(X_*^T). \end{aligned}$$

The proof of Theorem 9.13 shows first that  $(\lambda_*^t, -(\lambda_*^t \cdot f(X_*^t, \alpha_*^t) + h(X_*^t, \alpha_*^t)))$  is the first variation in  $x$  and  $t$  of the value function at the point  $(X_*^t, t)$  and concludes then that the first variation is in the superdifferential, that is

$$\left( \lambda_*^t, -(\lambda_*^t \cdot f(X_*^t, \alpha_*^t) + h(X_*^t, \alpha_*^t)) \right) \in D^+u(X_*^t, t).$$

Since the value function is a viscosity solution we conclude that

$$-(\lambda_*^t \cdot f(X_*^t, \alpha_*^t) + h(X_*^t, \alpha_*^t)) + \underbrace{H(\lambda_*^t, x)}_{\min_{\alpha \in A} (\lambda_*^t \cdot f(X_*^t, \alpha) + h(X_*^t, \alpha))} \geq 0$$

which means that  $\alpha_*$  satisfies the Pontryagin principle also in the case of non differentiable solutions to Hamilton-Jacobi equations.

### 9.1.6.2 Semiconcave Value Functions

There is an alternative and maybe more illustrative proof of the last theorem in a special setting: namely when the set of backward optimal paths  $\{(\bar{X}^t, t) : t < T\}$ , solving (9.29) and (9.47), may collide into a codimension one surface  $\Gamma$  in space-time  $\mathbb{R}^d \times [0, T]$ . Assume the value function is attained by precisely one path for  $(x, t) \in \mathbb{R}^d \times [0, T] - \Gamma$  and that the minimum is attained by precisely two paths at  $(x, t) \in \Gamma$ . Colliding backward paths (or characteristics)  $X$  in general lead to a discontinuity in the gradient of the value function,  $\lambda = u_x$ , on the surface of collision, which means that the surface is a shock wave for the multidimensional system of conservation laws

$$\partial_t \lambda^i(x, t) + \frac{d}{dx_i} H(\lambda(x, t), x) = 0 \quad (x, t) \in \mathbb{R}^d \times [0, T], \quad i = 1, \dots, d.$$

Denote the jump, for fixed  $t$ , of a function  $w$  at  $\Gamma$  by  $[w]$ . To have two colliding paths at a point on  $\Gamma$  requires that  $\lambda$  has a jump  $[\lambda] \neq 0$  there, since  $[\lambda] = 0$  yields only one path. The implicit function theorem shows that for fixed  $t$  any compact subset of the set  $\Gamma(t) \equiv \Gamma \cap (\mathbb{R}^d \times \{t\})$  is a  $\mathcal{C}^1$  surface: the surface  $\Gamma(t)$  is defined by the value functions,  $u^1$

and  $u^2$  for the two paths colliding on  $\Gamma$ , being equal on  $\Gamma$  and there are directions  $\hat{n} \in \mathbb{R}^d$  so that the Jacobian determinant  $\hat{n} \cdot \nabla(u^1 - u^2) = \hat{n} \cdot [\lambda] \neq 0$ . Therefore compact subsets of the surface  $\Gamma(t)$  has a well defined unit normal  $n$ . We assume that  $\Gamma(t)$  has a normal everywhere and we will prove that  $[\lambda] \cdot n \leq 0$ , which implies that  $u$  is semi-concave.

Two optimal backwards paths that collide on  $(x, t) \in \Gamma$  must depart in opposite direction away from  $\Gamma$ , that is  $n \cdot H_\lambda(\lambda_+, x) \geq 0$  and  $n \cdot H_\lambda(\lambda_-, x) \leq 0$ , see Figure 9.6, so that

$$0 \leq n \cdot [H_\lambda(\lambda, x)] = n \cdot \underbrace{\int_0^1 H_{\lambda\lambda}(\lambda_- + s[\lambda]) ds[\lambda]}_{=: \bar{H}_{\lambda\lambda} \leq 0}. \quad (9.19)$$

We know that  $u$  is continuous also around  $\Gamma$ , therefore the jump of the gradient,  $[u_x]$ , has to be parallel to the normal,  $n$ , of the surface  $\Gamma$ . Lemma 9.28 shows that  $[u_x] = [\lambda]$  and we conclude that this jump  $[\lambda]$  is parallel to  $n$  so that  $[\lambda] = [\lambda \cdot n]n$ , which combined with (9.19) shows that

$$0 \leq [\lambda \cdot n] \bar{H}_{\lambda\lambda} n \cdot n.$$

The  $\lambda$ -concavity of the Hamiltonian, see Exercise 9.19, implies that the matrix  $H_{\lambda\lambda}$  is negative semidefinite and consequently

$$\bar{H}_{\lambda\lambda} n \cdot n \leq 0, \quad (9.20)$$

which proves the claim  $[\lambda] \cdot n \leq 0$ , if we can exclude equality in (9.20). Equality in (9.20) means that  $\bar{H}_{\lambda\lambda} n = 0$  and implies  $H_\lambda(\lambda^+(t), x) = H_\lambda(\lambda^-(t), x)$  which is not compatible with two outgoing backward paths. Hence equality in (9.20) is ruled out. This derivation can be extended to several paths colliding into one point, see Exercise 9.18.

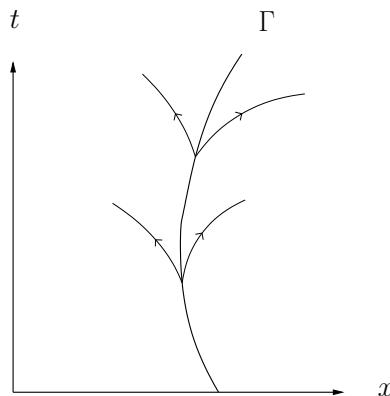


Figure 9.6: Optimal paths departing away from  $\Gamma$ .

**Exercise 9.18.**



**Exercise 9.19.** Show that the Hamiltonian

$$H(\lambda, x) := \min_{\alpha \in A} (\lambda \cdot f(x, \alpha) + h(x, \alpha))$$

is concave in the  $\lambda$ -variable, that is show that for each  $\lambda^1$  and  $\lambda^2$  in  $\mathbb{R}^d$  and for all  $s \in [0, 1]$  there holds

$$H(s\lambda^1 + (1-s)\lambda^2, x) \geq sH(\lambda^1, x) + (1-s)H(\lambda^2, x).$$

### 9.1.7 Maximum Norm Stability of Viscosity Solutions

An important aspect of the viscosity solution of the Hamilton-Jacobi-Bellman equation is its maximum norm stability with respect to maximum norm perturbations of the data, in this case the Hamiltonian and the initial data; that is the value function is stable with respect to perturbations of the flux  $f$  and cost functions  $h$  and  $g$ .

Assume first for simplicity that the optimal control is attained and that the value function is differentiable for two different optimal control problems with data  $f, h, g$  and the Hamiltonian  $H$ , respectively  $\bar{f}, \bar{h}, \bar{g}$  and Hamiltonian  $\bar{H}$ . The general case with only superdifferentiable value functions is studied afterwards. We have for the special case with the same initial data  $\bar{X}^0 = X^0$  and  $\bar{g} = g$

$$\begin{aligned} & \underbrace{\int_0^T \bar{h}(\bar{X}^t, \bar{\alpha}^t) dt + \bar{g}(\bar{X}^T)}_{\bar{u}(\bar{X}^0, 0)} - \underbrace{\int_0^T h(X^t, \alpha^t) dt + g(X^T)}_{u(X^0, 0)} \\ &= \int_0^T \bar{h}(\bar{X}^t, \bar{\alpha}^t) dt + u(\bar{X}^T, T) - \underbrace{u(X^0, 0)}_{u(\bar{X}^0, 0)} \\ &= \int_0^T \bar{h}(\bar{X}^t, \bar{\alpha}^t) dt + \int_0^T du(\bar{X}^t, t) \tag{9.21} \\ &= \int_0^T \underbrace{\partial_t u(\bar{X}^t, t)}_{=-H(\partial_x u(\bar{X}^t, t), \bar{X}^t)} + \underbrace{\partial_x u(\bar{X}^t, t) \cdot \bar{f}(\bar{X}^t, \bar{\alpha}^t) + \bar{h}(\bar{X}^t, \bar{\alpha}^t)}_{\geq \bar{H}(\partial_x u(\bar{X}^t, t), \bar{X}^t)} dt \\ &\geq \int_0^T (\bar{H} - H)(\partial_x u(\bar{X}^t, t), \bar{X}^t) dt. \end{aligned}$$

The more general case with  $\bar{g} \neq g$  yields the additional error term

$$(g - \bar{g})(\bar{X}^T)$$

to the right hand side in (9.21).

To find an upper bound, repeat the derivation above, replacing  $u$  along  $\bar{X}^t$  with  $\bar{u}$

along  $X^t$ , to obtain

$$\begin{aligned}
& \underbrace{\int_0^T h(X^t, \alpha^t) dt + g(X^T)}_{u(X^0, 0)} - \underbrace{\int_0^T \bar{h}(\bar{X}^t, \bar{\alpha}^t) dt + \bar{g}(\bar{X}^T)}_{\bar{u}(\bar{X}^0, 0)} \\
&= \int_0^T h(X^t, \alpha^t) dt + \bar{u}(X^T, T) - \underbrace{\bar{u}(\bar{X}^0, 0)}_{\bar{u}(X^0, 0)} \\
&= \int_0^T h(X^t, \alpha^t) dt + \int_0^T d\bar{u}(X^t, t) \\
&= \int_0^T \underbrace{\partial_t \bar{u}(X^t, t)}_{=-\bar{H}(\partial_x \bar{u}(X^t, t), X^t)} + \underbrace{\partial_x \bar{u}(X^t, t) \cdot f(X^t, \alpha^t) + h(X^t, \alpha^t)}_{\geq H(\partial_x \bar{u}(X^t, t), X^t)} dt \\
&\geq \int_0^T (H - \bar{H})(\partial_x \bar{u}(X^t, t), X^t) dt.
\end{aligned}$$

The two estimates above yields both an upper and a lower bound

$$\begin{aligned}
\int_0^T (H - \bar{H})(\partial_x \bar{u}(X^t, t), X^t) dt &\leq u(X_0, 0) - \bar{u}(X_0, 0) \\
&\leq \int_0^T (H - \bar{H})(\partial_x u(\bar{X}^t, t), \bar{X}^t) dt.
\end{aligned} \tag{9.22}$$

**Remark 9.20** (No minimizers). If there are no minimizers  $(\alpha, X)$  and  $(\bar{\alpha}, \bar{X})$ , then for every  $\varepsilon > 0$ , we can choose controls  $\alpha, \bar{\alpha}$  with corresponding states  $X, \bar{X}$  such that

$$E_{lhs} - \varepsilon \leq u(X_0, 0) - \bar{u}(X_0, 0) \leq E_{rhs} + \varepsilon$$

with  $E_{lhs}, E_{rhs}$  being the left and right hand sides of (9.22).

Solutions to Hamilton-Jacobi equations are in general not differentiable as we have seen in Example 9.9. Let us extend the derivation of (9.22) to a case when  $u$  is not differentiable. If  $u$  is a non differentiable semiconcave solution to a Hamilton-Jacobi equation, Definition 9.12 of the viscosity solution reduces to

$$\begin{aligned}
p_t + H(p_x, x) &= 0 \quad \text{for all } (p_t, p_x) \in Du(x, t) \text{ and all } t < T, x \in \mathbb{R}^d, \\
p_t + H(p_x, x) &\geq 0 \quad \text{for all } (p_t, p_x) \in D^+u(x, t) \text{ and all } t < T, x \in \mathbb{R}^d, \\
u(\cdot, T) &= g.
\end{aligned}$$

Consider now a point  $(x, t)$  where the value function is not differentiable. This means that in (9.21) we can for each  $t$  choose a point  $(p_t, p_x) \in D^+u(X^t, t)$  so that

$$\begin{aligned}
\int_0^T du(\bar{X}^t, t) + \int_0^T \bar{h}(\bar{X}^t, \bar{\alpha}^t) dt &= \int_0^T (p_t + p_x \cdot \bar{f}(\bar{X}^t, \bar{\alpha}^t) + \bar{h}(\bar{X}^t, \bar{\alpha}^t)) dt \\
&\geq \int_0^T (p_t + \bar{H}(p_x, \bar{X}^t)) dt \geq \int_0^T (-H + \bar{H})(p_x, \bar{X}^t) dt.
\end{aligned}$$

Note that the only difference compared to the differentiable case is the inequality instead of equality in the last step, which uses that optimal control problems have semi-concave viscosity solutions. The analogous formulation holds for  $\bar{u}$ . Consequently (9.22) holds for some  $(p_t, p_x) \in D^+u(\bar{X}^t, t)$  replacing  $(\partial_t u(\bar{X}^t, t), \partial_x u(\bar{X}^t, t))$  and some  $(\bar{p}_t, \bar{p}_x) \in D^+\bar{u}(X^t, t)$  replacing  $(\partial_t \bar{u}(X^t, t), \partial_x \bar{u}(X^t, t))$ .

The present analysis is in principle valid even when we replace  $\mathbb{R}^d$  to be an infinite dimensional Hilbert space for optimal control of partial differential equations, although existence and semiconcavity of solutions is not derived in full generality, see [San08]

## 9.2 Numerical Approximation of ODE Constrained Minimization

We consider numerical approximations with the time steps

$$t_n = \frac{n}{N}T, \quad n = 0, 1, 2, \dots, N.$$

The most basic approximation is based on the minimization

$$\min_{\bar{\alpha} \in B^N} \left( g(\bar{X}_N) + \sum_{n=0}^{N-1} h(\bar{X}_n, \bar{\alpha}_n) \Delta t \right), \quad (9.23)$$

where  $\Delta t = t_{n+1} - t_n$ ,  $\bar{X}_0 = X_0$  and  $\bar{X}_n \equiv \bar{X}(t_n)$ , for  $1 \leq n \leq N$ , satisfy the *forward Euler* constraint

$$\bar{X}_{n+1} = \bar{X}_n + \Delta t f(\bar{X}_n, \bar{\alpha}_n). \quad (9.24)$$

The existence of at least one minimum of (9.23) is clear since it is a minimization of a continuous function in the compact set  $B^N$ . The Lagrange principle can be used to solve such a constrained minimization problem. We will focus on a variant of this method based on the *discrete Pontryagin principle* where the control is eliminated

$$\begin{aligned} \bar{X}_{n+1} &= \bar{X}_n + \Delta t H_\lambda(\bar{\lambda}_{n+1}, \bar{X}_n), & \bar{X}_0 &= X_0, \\ \bar{\lambda}_n &= \bar{\lambda}_{n+1} + \Delta t H_x(\bar{\lambda}_{n+1}, \bar{X}_n), & \bar{\lambda}_N &= g_x(\bar{X}_N), \end{aligned} \quad (9.25)$$

called the symplectic Euler method for the Hamiltonian system (9.13), cf. [HLW02].

A natural question is in what sense the discrete problem (9.25) is an approximation to the continuous optimal control problem (9.13). In this section we show that the value function of the discrete problem approximates the continuous value function, using the theory of viscosity solutions to Hamilton-Jacobi equations to construct and analyse regularized Hamiltonians.

Our analysis is a kind of backward error analysis. The standard backward error analysis for Hamiltonian systems uses an analytic Hamiltonian and shows that symplectic one step schemes generate approximate paths that solve a modified Hamiltonian system, with the perturbed Hamiltonian given by a series expansion cf. [HLW02]. Our backward error analysis is different and more related to the standard finite element analysis. We

first extend the approximate Euler solution to a continuous piecewise linear function in time and define a *discrete value function*,  $\bar{u} : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}$ . This value function satisfies a perturbed Hamilton-Jacobi partial differential equation, with a small residual error. A special case of our analysis shows that if the optimal  $\alpha$  in (9.5) is a differentiable function of  $x$  and  $\lambda$  and if the optimal backward paths,  $\bar{X}(s)$  for  $s < T$ , do not collide, more about this later, the discrete value functions,  $\bar{u}$ , for the Pontryagin method (9.25) satisfies a Hamilton-Jacobi equation:

$$\bar{u}_t + H(\bar{u}_x, \cdot) = \mathcal{O}(\Delta t), \quad \text{as } \Delta t \rightarrow 0+, \quad (9.26)$$

where

$$\bar{u}(x, t_m) \equiv \min_{\bar{\alpha} \in B^N} \left( g(\bar{X}_N) + \sum_{n=m}^{N-1} h(\bar{X}_n, \bar{\alpha}_n) \Delta t \right) \quad (9.27)$$

for solutions  $\bar{X}$  to with  $\bar{X}(t_m) \equiv \bar{X}_m = x$ . The minimum in (9.27) is taken over the solutions to the discrete Pontryagin principle (9.25). The maximum norm stability of Hamilton–Jacobi PDE solutions and a comparison between the two equations (9.4) and (9.26) show that

$$\mathcal{O} \|u - \bar{u}\|_{\mathcal{C}} = \mathcal{O}(\Delta t). \quad (9.28)$$

However, in general the optimal controls  $\bar{\alpha}$  and  $\alpha$  in (9.24) and (9.1) are discontinuous functions of  $x$ , and  $\bar{\lambda}$  or  $u_x$ , respectively, and the backward paths *do* collide. There are two different reasons for discontinuous controls:

- The Hamiltonian is in general only Lipschitz continuous, even if  $f$  and  $h$  are smooth.
- The optimal backward paths may collide.

The standard error analysis for ordinary differential equations is directly applicable to control problems when the time derivative of the control function is integrable. But general control problems with discontinuous controls require alternative analysis, which will be in two steps. The *first step* in our error analysis is to construct regularizations of the functions  $f$  and  $h$ , based on (9.14) applied to a  $\mathcal{C}^2(\mathbb{R}^d \times \mathbb{R}^d)$  approximate Hamiltonian  $H^\delta$  which is  $\lambda$ -concave and satisfies

$$\|H^\delta - H\|_{\mathcal{C}} = \mathcal{O}(\delta), \quad \text{as } \delta \rightarrow 0^+,$$

and to introduce the regularized paths

$$\begin{aligned} \bar{X}_{n+1} &= \bar{X}_n + \Delta t H_\lambda^\delta(\bar{\lambda}_{n+1}, \bar{X}_n), & \bar{X}_0 &= X_0, \\ \bar{\lambda}_n &= \bar{\lambda}_{n+1} + \Delta t H_x^\delta(\bar{\lambda}_{n+1}, \bar{X}_n), & \bar{\lambda}_N &= g_x(\bar{X}_N). \end{aligned} \quad (9.29)$$

We will sometimes use the notation  $f^\delta \equiv H_\lambda^\delta$  and  $h^\delta \equiv H^\delta - \lambda H_\lambda^\delta$ .

The *second step* is to estimate the residual of the discrete value function in the Hamilton-Jacobi-Bellman equation (9.4). The maximum norm stability of viscosity solutions and the residual estimate imply then an estimate for the error in the value

function. An approximation of the form (9.29) may be viewed as a general symplectic one step method for the Hamiltonian system (9.13), see Section 9.2.7.

There is a second reason to use Hamiltonians with smooth flux: in practice the nonlinear boundary value problem (9.29) has to be solved by iterations. If the flux is not continuous it seems difficult to construct a convergent iterative method, in any case iterations perform better with smoother solutions. When the Hamiltonian can be formed explicitly, the Pontryagin based method has the advantage that the Newton method can be applied to solve the discrete nonlinear Hamiltonian system with a sparse Jacobian.

If the optimal discrete backward paths  $\bar{X}(t)$  in (9.29) collide on a codimension one surface  $\Gamma$  in  $\mathbb{R}^d \times [0, T]$ , the dual variable  $\bar{\lambda} = \bar{u}_x$  may have a discontinuity at  $\Gamma$ , as a function of  $x$ . Theorems 9.27 and ?? prove, for  $\bar{u}$  based on the Pontryagin method, that in the viscosity solution sense

$$\bar{u}_t + H(\bar{u}_x, \cdot) = \mathcal{O}(\Delta t + \delta + \frac{(\Delta t)^2}{\delta}), \quad (9.30)$$

where the discrete value function,  $\bar{u}$ , in (9.27) has been modified to

$$\bar{u}(x, t_m) = \min_{\bar{X}_m=x} \left( g(\bar{X}_N) + \sum_{n=m}^{N-1} h^\delta(\bar{X}_n, \bar{\lambda}_{n+1}) \Delta t \right). \quad (9.31)$$

The regularizations make the right hand side in (9.30) a Lipschitz continuous function of  $(\bar{\lambda}(t), \bar{X}(t), t)$ , bounded by  $C(\Delta t + \delta + \frac{(\Delta t)^2}{\delta})$  where  $C$  depends only on the Lipschitz constants of  $f$ ,  $h$  and  $\bar{\lambda}$ . Therefore the maximum norm stability can be used to prove  $\|u - \bar{u}\|_C = \mathcal{O}(\Delta t)$ , for  $\delta = \Delta t$ . Without the regularization, the corresponding error term to in (9.30) is not well defined, even if  $\bar{u}_x$  is smooth. A similar proof applies to the minimization method for smooth Hamiltonians, see [San08]. It is important to note that for non smooth control the solution paths  $\bar{X}$  may not converge although the value function converges as  $\Delta t$  and  $\delta$  tend to zero. Therefore our backward error analysis uses consistency with the Hamilton-Jacobi partial differential equation and not with the Hamiltonian system. Convergence of the approximate path  $(\bar{X}, \bar{\lambda})$  typically requires Lipschitz continuous flux  $(H_\lambda, H_x)$ , which we do not assume in this work.

### 9.2.1 Optimization Examples

We give some examples when the Hamiltonian,  $H$ , is not a differentiable function, and difficulties associated with this.

**Example 9.21.** Let  $B = \{-1, 1\}$ ,  $f = \alpha$ ,  $h = x^2/2$  and  $g = 0$ . Here the continuous minimization problem (9.3) has no minimizer among the measurable functions. A solution in discrete time using a nonregularized Pontryagin method or discrete dynamic programming will behave as in Figure 9.7. First the solution approaches the time axis, and then it oscillates back and forth. As  $\Delta t$  becomes smaller these oscillations do so as well. The infimum for the continuous problem corresponds to a solution  $X(t)$  that approaches the time-axis, and then remains on it. However, this corresponds to  $\alpha = 0$ ,

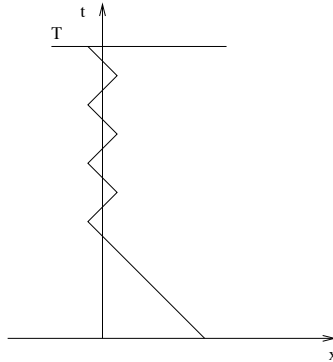


Figure 9.7: Example 9.21 where the continuous problem has no minimizer among the measurable functions.

which is not in  $B$ , and hence the infimum is not attained. A cure to always have an attained minimizing path for the continuous problem is to use controls which are Young measures, see [You69] and [Ped99]. We note that the Hamiltonian,  $H(\lambda, x) = -|\lambda| + x^2/2$ , in this example is not differentiable.

**Example 9.22.** Let  $B = [-1, 1]$ ,  $f = \alpha$ ,  $h = x^2/2$  and  $g = 0$ , which is similar to the previous example but now the set of admissible controls,  $B$ , has been changed slightly. Since  $0 \in B$ , the infimum in (9.3) is now obtained. However, the Hamiltonian remains unchanged compared to the previous example, and a solution to the discrete Pontryagin principle would still be oscillating as in Figure 9.7.

**Example 9.23.** Let  $B = [-1, 1]$ ,  $f = \alpha$ ,  $h = 0$  and  $g = x^2$ . The Hamiltonian is nondifferentiable:  $H = -|\lambda|$ . If  $T = 1$  there are infinitely many solutions to the continuous minimization, the discrete minimization and the unregularized discrete Pontryagin principle, when  $X_0 \in (-1, 1)$ , as depicted in Figure 9.8.

The problems occurring in the previous examples are all cured by regularizing the Hamiltonian and using the scheme (9.29). That is, the solution to (9.29) in the first two examples is a smooth curve that obtains an increasingly sharp kink near the time-axis as the regularizing parameter,  $\delta$ , decreases, see Figure 9.9. In the last of the previous examples we, in contrast to the other methods, obtain a unique solution to (9.29).

Another problem that has not to do with nondifferentiability of the Hamiltonian is shown in the following example:

**Example 9.24.** Let  $B = [-1, 1]$ ,  $f = \alpha$ ,  $h = 0$  and  $g = -|x|$ . Although  $H$  is discontinuous here, this is not what causes the problem. The problem is that optimal paths collide backwards, see Figure 9.10. When  $X_0 = 0$  there are two solutions, one going to the left, and one to the right. The left solution has  $\lambda = 1$  and the right solution has  $\lambda = -1$ , so on the time-axis  $\lambda$  is discontinuous. For these values of  $\lambda$ , the Hamiltonian is differentiable, therefore the nonsmoothness of the Hamiltonian is not the issue here. It is rather the global properties of the problem that play a role. This problem is difficult to

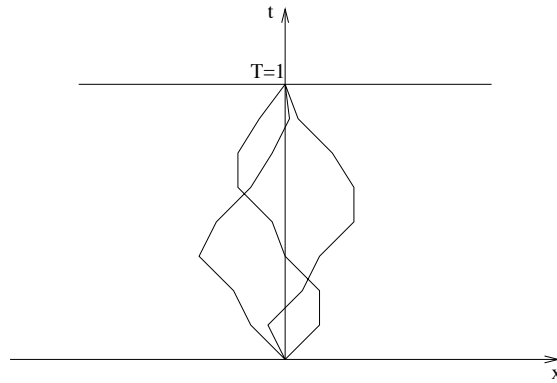


Figure 9.8: Example 9.23 with  $g(x) = x^2$  gives infinitely many minimizing paths through the same starting point.

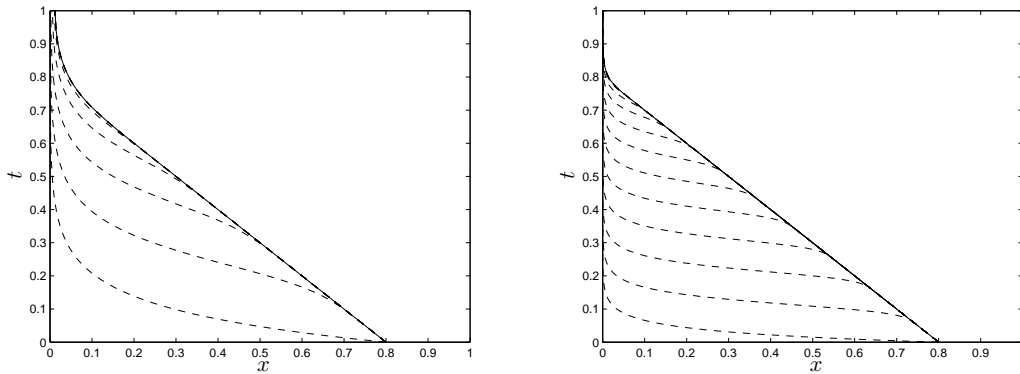


Figure 9.9: Solution of the discrete optimization problem (9.29) in Example 9.21 and 9.22 for  $\delta = \Delta t = 1/N$ ,  $X_0 = 0.8$  and  $H_\lambda^\delta(\lambda, x) = -\tanh(\lambda/\delta)$ , using the Newton method. To the left,  $N = 100$ , and to the right,  $N = 1000$ . The dashed lines shows the solution after each Newton iteration.

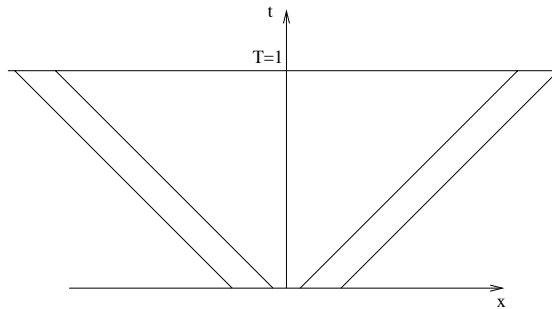


Figure 9.10: Solution of the optimization problem in Example 9.24, where  $g(x) = -|x|$ ,  $f = \alpha$ ,  $h = 0$  and  $B = [-1, 1]$ , for four different starting points.

regularize, and it will not be done here. However, we still can show convergence of the scheme (9.29). This is done in Section ??.

When using (9.29) to solve the minimization problem (9.3) it is assumed that the Hamiltonian is exactly known. Is this an unrealistic assumption in practice? In the following two examples we indicate that there exist interesting examples where we know the Hamiltonian. The first has to do with volatility estimation in finance, and the latter with optimization of an electric contact.

### 9.2.1.1 Implied Volatility

Black-Scholes equation for pricing general options uses the volatility of the underlying asset. This parameter, however, is difficult to estimate. One way of estimation is to use measured market values of options on the considered asset for standard European contracts. This way of implicitly determining the volatility is called implied volatility. In the simplest setting, the formula<sup>2</sup> for the option price based on constant interest rate and volatility is used. Then the result typically gives different values of the volatility for different stock price – instead of obtaining a constant volatility, the implied volatility becomes a strictly convex function of the stock price called the volatility smile. Below we shall fit a model allowing the volatility to be a general function to observed option prices. That requires solution of a partial differential equation, since an explicit formula is not available. Another ingredient in our reconstruction is to use the so called Dupire equation for standard European put and call option prices as a function of the strike price and strike time. Using an equation of the option value as a function of the strike price and strike time, for given stock price, is computational more efficient, since the option data is for different strike price and strike times, with fixed stock price. To use the standard Black-Scholes equation for the option value as a function of the stock price

<sup>2</sup>the option price formula is  $C(s, t; K, T) = s\Phi(d_1) - Ke^{-r(T-t)}\Phi(d_2)$ , where  $d_1 := (\ln(s/K) + (r + \sigma^2/2)(T-t))/(\sigma(T-t)^{1/2})$ ,  $d_2 := d_1 - \sigma(T-t)^{1/2}$  and  $\Phi$  is the standard normal cumulative distribution function.



would require to solve different equations for each data point, which is also possible but more computationally expensive.

We assume that the financial asset obeys the following Ito stochastic differential equation,

$$dS(t) = \mu S(t)dt + \sigma(t, S(t))S(t)dW(t), \quad (9.32)$$

where  $S(t)$  is the price of the asset at time  $t$ ,  $\mu$  is a drift term,  $\sigma$  is the volatility and  $W : \mathbb{R}_+ \rightarrow \mathbb{R}$  is the Wiener process. If the volatility is a sufficiently regular function of  $S$ ,  $t$ , the strike level  $K$  and the maturity date  $T$ , the Dupire equation holds for the option price  $C(T, K)$  as a function of  $T$  and  $K$ , with the present time  $t = 0$  and stock price  $S(0) = S$  fixed,

$$\begin{aligned} C_T - \tilde{\sigma} C_{KK} &= 0, \quad T \in (0, \infty), K > 0, \\ C(0, K) &= \max\{S - K, 0\} \quad K > 0, \end{aligned} \quad (9.33)$$

where

$$\tilde{\sigma}(T, K) \equiv \frac{\sigma^2(T, K)K^2}{2}.$$

Here the contract is an european call option with payoff function  $\max\{S(T) - K, 0\}$ . We have for simplicity assumed the bank rate to be zero. A derivation of Dupire's equation (9.33) is presented in Example 9.25 in the special setting  $r = 0$ ; the general case is studied in [Dup94].

The optimization problem now consists of finding  $\sigma(T, K)$  such that

$$\int_0^{\hat{T}} \int_{\mathbb{R}_+} (C - \hat{C})^2(T, K)w(T, K)dKdT \quad (9.34)$$

is minimized, where  $\hat{C}$  are the measured market values on option prices for different strike prices and strike times and  $w$  is a non negative weight function. In practice,  $\hat{C}$  is not known everywhere, but for the sake of simplicity, we assume it is and set  $w \equiv 1$ , that is there exists a future time  $\hat{T}$  such that  $\hat{C}$  is defined in  $\mathbb{R}_+ \times [0, \hat{T}]$ . If the geometric Brownian motion would be a perfect model for the evolution of the price of the asset, the function  $\sigma(T, K)$  would be constant, but as this is not the case, the  $\sigma$  that minimizes (9.34) (if a minimizer exists) varies with  $T$  and  $K$ .

It is possible to use (9.13) and (9.25) to perform the minimization of (9.34) over the solutions to a finite difference discretization of (9.33)

$$\begin{aligned} \min_{\tilde{\sigma}} \int_0^{\hat{T}} \Delta K \sum_i (C - \hat{C})_i^2 dT \\ \text{subject to } \frac{\partial C_i(T)}{\partial T} &= \tilde{\sigma} D^2 C_i(T), \\ C_i(0) &= \max(S - i\Delta K, 0), \end{aligned} \quad (9.35)$$

where we now let  $C_i(T) \approx C(T, i\Delta K)$  denote the discretized prize function, for strike time  $T$  and strike price  $i\Delta K$ , and  $D^2$  is the standard three point difference approximation

of the second order partial derivative in  $K$ , that is  $(D^2C)_i = (C_{i+1} - 2C_i + C_{i-1})/\Delta K^2$ . In order to have a finite dimensional problem we restrict to a compact interval  $(0, M\Delta K)$  in  $K$  with the boundary conditions

$$C_0 = S, \quad C_M = 0.$$

This formulation will be exactly the same as in (9.13) if  $\Delta K = 1$ , and otherwise it requires to use a new scalar product  $(x, y) := \Delta K \sum_i x_i y_i$  and let the partial derivative  $\partial_\lambda$  be replaced by the following Gateaux derivative,  $H_\lambda$ ,

$$\lim_{\epsilon \rightarrow 0} \epsilon^{-1} (H(\lambda + \epsilon v, C) - H(\lambda, C)) =: (H_\lambda(\lambda, C), v),$$

and similarly for  $\partial_C$ ; so that the partial derivative is a factor of  $\Delta K$  smaller than the Gateaux derivative. This complication with using  $\Delta K \neq 1$  is introduced in order to have a consistent formulation with the infinite dimensional case, where a partial derivative of a functional becomes zero but the Gateaux derivative is nonzero and meaningful, see the next example. The reader may avoid this by considering  $\Delta K = 1$ .

The Hamiltonian for this problem is

$$\begin{aligned} H(\lambda, C) &= \Delta K \min_{\tilde{\sigma}} \sum_{i=1}^{M-1} \left( \lambda_i \tilde{\sigma}_i (D^2C)_i + (C - \hat{C})_i^2 \right) \\ &= \Delta K \sum_{i=1}^{M-1} \left( \min_{\tilde{\sigma}_i} \lambda_i \tilde{\sigma}_i (D^2C)_i + (C - \hat{C})_i^2 \right) \end{aligned}$$

where  $\lambda$  is the adjoint associated to the constraint (9.35). We have used that the components of the flux,  $f$ , in this problem is  $\tilde{\sigma}_i (D^2C)_i$ , that the running cost,  $h$ , is  $\Delta K \sum_i (C - \hat{C})_i^2$ , and further that each  $\tilde{\sigma}_i$  minimizes  $\lambda_i \tilde{\sigma}_i (D^2C)_i$  separately, so that the minimum can be moved inside the sum. If we make the simplifying assumption that  $0 \leq \sigma_- \leq \tilde{\sigma} \leq \sigma_+ < \infty$  we may introduce a function  $s : \mathbb{R} \rightarrow \mathbb{R}$  as

$$s(y) \equiv \min_{\tilde{\sigma}} y \tilde{\sigma} = \begin{cases} y\sigma_-, & y > 0 \\ y\sigma_+, & y < 0. \end{cases}$$

Using  $s$ , it is possible to write the Hamiltonian as

$$H(\lambda, C) = \Delta K \sum_{i=1}^{M-1} \left( s(\lambda_i (D^2C)_i) + (C - \hat{C})_i^2 \right).$$

Since  $s$  is nondifferentiable, so is  $H$ . However,  $s$  may easily be regularized, and it is possible to obtain the regularization in closed form, e.g. as in Example 1. Using a regularized version  $s_\delta$  of  $s$ , the regularized Hamiltonian becomes

$$H^\delta(\lambda, C) = \Delta K \sum_{i=1}^{M-1} \left( s_\delta(\lambda_i (D^2C)_i) + (C - \hat{C})_i^2 \right),$$

which using Gateaux derivatives gives the Hamiltonian system

$$\begin{aligned}\frac{\partial C_i(T)}{\partial T} &= s'_\delta(\lambda_i(D^2C)_i)D^2C_i(T), \quad C_0 = S \quad C_M = 0, \\ -\frac{\partial \lambda_i(T)}{\partial T} &= D^2\left(s'_\delta(\lambda_i(D^2C)_i)\lambda\right) + 2(C - \hat{C})_i, \\ \lambda_0 &= \lambda_M = 0,\end{aligned}\tag{9.36}$$

with data

$$C_i(0) = \max(S - i\Delta K, 0), \quad \lambda(\hat{T}) = 0.$$

The corresponding Hamilton-Jacobi equation for the value function

$$u(C, \tau) = \int_\tau^{\hat{T}} \sum_{i=1}^{M-1} (C - \hat{C})_i^2 \Delta K dT$$

is

$$\begin{aligned}u_T + H(u_C, C) &= 0, \quad T < \hat{T}, \\ u(\hat{T}, \cdot) &= 0,\end{aligned}$$

where  $u_C$  is the Gateaux derivative with respect to  $C$  in the scalar product  $(x, y) \equiv \Delta K \sum_i x_i, y_i$ . With this scalar product the Hamiltonian system (9.36) takes the form

$$\begin{aligned}(C_T, v) &= (H_\lambda^\delta, v), \quad \forall v \in \mathbb{R}^{M-1} \\ (\lambda_T, v) &= -(H_C^\delta, v), \quad \forall v \in \mathbb{R}^{M-1}\end{aligned}$$

where  $H_\lambda^\delta$  and  $H_C^\delta$  are the Gateaux derivatives.

A choice of the regularization parameter  $\delta$ , depending also on data error, can be obtained e.g. by the discrepancy principle, cf. [Vog02], [EHN96]. The Newton method described in Section 3 works well to solve the discrete equations for  $d = 10$ . The results of one trial volatility estimation is given in Figure 9.11.

**Example 9.25** (Derivation of Dupire's equation). The Black-Scholes equation for a general volatility function and interest  $r = 0$  is

$$\begin{aligned}\partial_t f + \frac{\sigma^2(s, t)s^2}{2} \partial_{ss} f &= 0 \quad t < T \\ f(s, T) &= \max(K - s, 0)\end{aligned}$$

which defines the option value  $f(s, t; K, T)$ . The goal is now to find the equation for  $f$  as a function of  $K$  and  $T$ . We know from the Kolmogorov backward equation that  $f(s, t; K, T) = \mathbb{E}[\max(K - S_T, 0) \mid S_t = s]$ , where  $dS_t = \sigma(S_t, t)S_t dW_t$ . The Kolmogorov forward equation shows that  $f(s, t; K, T) = \int_{\mathbb{R}} \max(K - y, 0)p(y, T; s, t)dy$  where

$$\begin{aligned}\partial_T p - \partial_{yy} \left( \frac{\sigma^2(y, T)y^2}{2} p \right) &= 0 \quad T > t \\ p(y, t; s, t) &= \delta(y - s).\end{aligned}$$

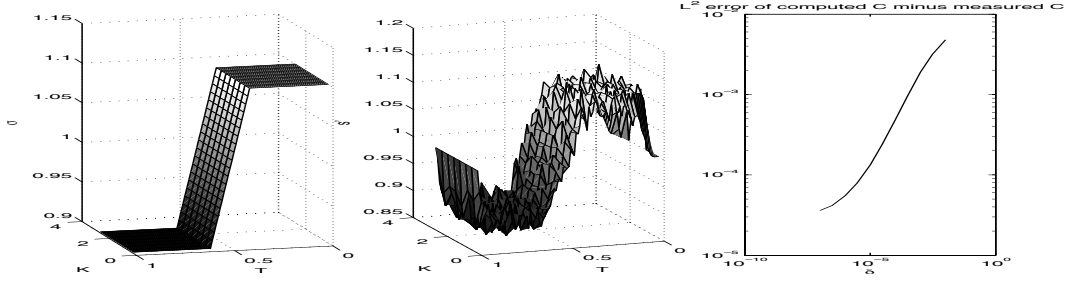


Figure 9.11: Results of a computer experiment where the volatility  $\sigma$  in the picture to the left is used to obtain the “measured”  $\hat{C}$ . Uniform noise of amplitude  $10^{-4}$  is also added to  $\hat{C}$ . The error  $\|C - \hat{C}\|_{L^2}$  is plotted versus  $\delta$  in the picture to the right. In the middle picture the approximate volatility,  $s'_\delta$  is shown for the value of  $\delta$  ( $= 3 \cdot 10^{-6}$ ) that minimizes  $\|s'_\delta - \sigma\|_{L^2}$ . In this experiment,  $M = 9$  and  $N = 100$ .

We observe that  $\partial_{KK}f(s, t; K, T) = \int_{\mathbb{R}} \delta(K - y)p(y, T; s, t)dy = p(K, T; s, t)$  and consequently

$$\partial_T \partial_{KK}f(s, t; K, T) - \partial_{KK} \left( \frac{\sigma^2(K, T)K^2}{2} \partial_{KK}f(s, t; K, T) \right) = 0 \quad T > t,$$

can be integrated to obtain

$$\partial_T f(s, t; K, T) - \left( \frac{\sigma^2(K, T)K^2}{2} \partial_{KK}f(s, t; K, T) \right) = C_1 + C_2 K \quad T > t.$$

The boundary condition  $\partial_{KK}f \rightarrow 0$  as  $K \rightarrow \infty$  and  $\partial_T f \rightarrow 0$  as  $T \rightarrow \infty$  concludes that  $C_1 = C_2 = 0$ .

### 9.2.1.2 Topology Optimization of Electric Conduction

The problem is to place a given amount of conducting material in a given domain  $\Omega \subset \mathbb{R}^d$  in order to minimize the power loss for a given surface current  $q$ , satisfying  $\int_{\partial\Omega} q ds = 0$ : let  $\eta \in \mathbb{R}$  be a given constant, associated to the given amount of material, and find an optimal conduction distribution  $\sigma : \Omega \rightarrow \{\sigma_-, \sigma_+\}$ , where  $\sigma_{\pm} > 0$ , such that

$$\begin{aligned} \operatorname{div}(\sigma \nabla \varphi(x)) &= 0, \quad x \in \Omega, & \sigma \frac{\partial \varphi}{\partial n} \Big|_{\partial\Omega} &= q \\ \min_{\sigma} \left( \int_{\partial\Omega} q \varphi ds + \eta \int_{\Omega} \sigma dx \right), \end{aligned} \tag{9.37}$$

where  $\partial/\partial n$  denotes the normal derivative and  $ds$  is the surface measure on  $\partial\Omega$ . Note that (9.37) implies that the power loss satisfies

$$\begin{aligned} \int_{\partial\Omega} q \varphi ds &= - \int_{\Omega} \operatorname{div}(\sigma \nabla \varphi) \varphi dx + \int_{\partial\Omega} \sigma \frac{\partial \varphi}{\partial n} \varphi ds \\ &= \int_{\Omega} \sigma \nabla \varphi \cdot \nabla \varphi dx. \end{aligned}$$

The Lagrangian takes the form

$$\int_{\partial\Omega} q(\varphi + \lambda) \, ds + \int_{\Omega} \sigma \underbrace{(\eta - \nabla\varphi \cdot \nabla\lambda)}_v \, dx$$

and the Hamiltonian becomes

$$H(\lambda, \varphi) = \min_{\sigma} \int_{\Omega} \sigma v \, dx + \int_{\partial\Omega} q(\varphi + \lambda) \, ds = \int_{\Omega} \underbrace{\min_{\sigma} \sigma v}_{s(v)} \, dx + \int_{\partial\Omega} q(\varphi + \lambda) \, ds$$

with the regularization

$$H^{\delta}(\lambda, \varphi) = \int_{\Omega} s_{\delta}(\eta - \nabla\varphi \cdot \nabla\lambda) \, dx + \int_{\partial\Omega} q(\varphi + \lambda) \, ds,$$

depending on the concave regularization  $s_{\delta} \in \mathcal{C}^2(\mathbb{R})$  as in Section 9.2.1.1. The value function

$$u(\varphi, \tau) = \int_{\tau}^T \left( \int_{\partial\Omega} q\varphi \, ds + \eta \int_{\Omega} \sigma \, dx \right) dt$$

for the parabolic variant of (9.37), that is

$$\varphi_t = \operatorname{div}(\sigma \nabla\varphi(x)),$$

yields the infinite dimensional Hamilton-Jacobi equation

$$\partial_t u + H(\partial_{\varphi} u, \varphi) = 0 \quad t < T, \quad u(\cdot, T) = 0,$$

using the Gateaux derivative  $\partial_{\varphi} u = \lambda$  of the functional  $u(\varphi, t)$  in  $L^2(\Omega)$ . The regularized Hamiltonian generates the following parabolic Hamiltonian system for  $\varphi$  and  $\lambda$

$$\begin{aligned} \int_{\Omega} (\partial_t \varphi w + s'(\eta - \nabla\varphi \cdot \nabla\lambda) \nabla\varphi \cdot \nabla w) \, dx &= \int_{\partial\Omega} q w \, ds \\ \int_{\Omega} (-\partial_t \lambda v + s'(\eta - \nabla\varphi \cdot \nabla\lambda) \nabla\lambda \cdot \nabla v) \, dx &= \int_{\partial\Omega} q v \, ds \end{aligned}$$

for all test functions  $v, w \in V \equiv \{v \in H^1(\Omega) \mid \int_{\Omega} v \, dx = 0\}$ . Time independent solutions satisfy  $\lambda = \varphi$  by symmetry. Therefore the electric potential satisfies the nonlinear elliptic partial differential equation

$$\operatorname{div}(s'_{\delta}(\eta - |\nabla\varphi|^2) \nabla\varphi(x)) = 0 \quad x \in \Omega, \quad s'_{\delta} \frac{\partial\varphi}{\partial n} \Big|_{\partial\Omega} = q, \quad (9.38)$$

which can be formulated as the convex minimization problem:  $\varphi \in V$  is the unique minimizer (up to a constant) of

$$- \left( \int_{\Omega} s_{\delta}(\eta - |\nabla\varphi(x)|^2) \, dx + 2 \int_{\partial\Omega} q\varphi \, ds \right). \quad (9.39)$$

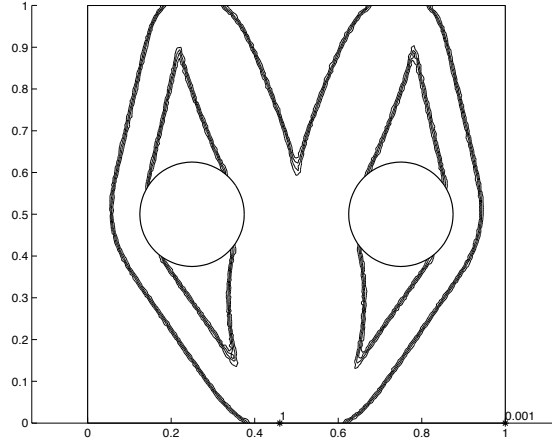


Figure 9.12: Contour plot of  $s'_\delta$  as an approximation of the conductivity  $\sigma$ . As seen,  $\Omega$  is in this example a square with two circles cut out. Electrical current enters  $\Omega$  at two positions on the top of the square and leaves at one position on the bottom. The contours represent the levels 0.2, 0.4, 0.6 and 0.8. A piecewise linear FEM was used with 31440 elements, maximum element diameter 0.01,  $\sigma_- = 0.001$ ,  $\sigma_+ = 1$ ,  $\eta = 0.15$  and  $\delta = 10^{-5}$ .

In [CSS08] we study convergence of

$$\lim_{T \rightarrow \infty} \frac{u(\varphi, t) - \bar{u}(\varphi, t)}{T},$$

where  $\bar{u}$  is the value function associated to finite element approximations of the minimization (9.39).

The Newton method in Section 3 works well to solve the finite element version of (9.38) by successively decreasing  $\delta$ , also for large  $d$ , see [CSS08], where also the corresponding inverse problem to use measured approximations of  $\varphi$  to determine the domain where  $\sigma = \sigma_-$  and  $\sigma = \sigma_+$  is studied. A numerical solution to (9.38) can be seen in Figure 9.12.

In this paper we use the standard Euclidean norm in  $\mathbb{R}^d$  to measure  $X$  and  $\lambda$ . Optimal control of partial differential equations with  $X$  and  $\lambda$  belonging to infinite dimensional function spaces requires a choice of an appropriate norm. In [San08] the analysis here is extended to optimal control of some parabolic partial differential equations, by replacing the Euclidean  $\mathbb{R}^d$  norm with the  $H_0^1$  Sobolev norm, using also that the theory of viscosity solutions remains valid with this replacement.

## 9.2.2 Solution of the Discrete Problem

We assume in the theorems that the Pontryagin minimization (9.29) has been solved exactly. In practice (9.29) can only be solved approximately by iterations. The simplest iteration method to solve the boundary value problem (9.29) is the shooting method:

start with an initial guess of  $\bar{\lambda}[0]$  and compute, for all time steps  $n$ , the iterates

$$\begin{aligned}\bar{X}_{n+1} &= \bar{X}_n + \Delta t H_\lambda^\delta(\bar{\lambda}_{n+1}[i], \bar{X}_n), \quad n = 0, \dots, N-1, \quad \bar{X}_0 = X_0 \\ \bar{\lambda}_n[i+1] &= \bar{\lambda}_{n+1}[i] + \Delta t H_x^\delta(\bar{\lambda}_{n+1}[i], \bar{X}_n), \quad n = N-1, \dots, 0, \quad \bar{\lambda}_N = g_x(\bar{X}_N).\end{aligned}\tag{9.40}$$

An alternative method, better suited for many boundary value problems, is to use Newton iterations for the nonlinear system  $F(\bar{X}, \bar{\lambda}) = 0$  where  $F : \mathbb{R}^{Nd} \times \mathbb{R}^{Nd} \rightarrow \mathbb{R}^{2Nd}$  and

$$\begin{aligned}F(\bar{X}, \bar{\lambda})_{2n} &= \bar{X}_{n+1} - \bar{X}_n - \Delta t H_\lambda^\delta(\bar{\lambda}_{n+1}, \bar{X}_n), \\ F(\bar{X}, \bar{\lambda})_{2n+1} &= \bar{\lambda}_n - \bar{\lambda}_{n+1} - \Delta t H_x^\delta(\bar{\lambda}_{n+1}, \bar{X}_n).\end{aligned}\tag{9.41}$$

An advantage with the Pontryagin based method (9.41) is that the Jacobian of  $F$  can be calculated explicitly and it is sparse. The Newton method can be used to solve the volatility and topology optimization examples in Section 2, where the parameter  $\delta$  is successively decreasing as the nonlinear equation (9.41) is solved more accurately.

Let us use dynamic programming to show that the system (9.29) has a solution in the case that  $\bar{\lambda}$  is a Lipschitz continuous function of  $(x, t)$ , with Lipschitz norm independent of  $\Delta t$ , and  $\delta > C\Delta t$ . One step

$$x = y + \Delta t H_\lambda^\delta(\lambda(x), y)\tag{9.42}$$

for fixed  $y \in \mathbb{R}^d$  has a solution  $x(y)$  since the iterations

$$x[i+1] = y + \Delta t H_\lambda^\delta(\lambda(x[i]), y)$$

yield a contraction for the error  $e[i] = x[i+m] - x[i]$

$$e[i+1] = \Delta t \left( H_\lambda^\delta(\lambda(x[i+m]), y) - H_\lambda^\delta(\lambda(x[i]), y) \right) = \Delta t \overline{H_{\lambda\lambda}^\delta} \lambda_x e[i].$$

Conversely, for all  $x \in \mathbb{R}^d$  equation (9.42) has a solution  $y(x)$  for each step since the iterations

$$y[i+1] = x - \Delta t H_\lambda^\delta(\lambda(x), y[i])$$

generate a contraction for the error. The dynamic programming principle then shows that there are unique paths through all points  $\bar{X}_{n+1}$  leading to all  $\bar{X}_n$  for all  $n$ .

**Example 9.26.** In Example 9.21 and 9.22 the problem was to minimize

$$\min_{\alpha \in B} \int_0^T \frac{X(t)^2}{2} dt,$$

given the dynamics

$$X'(t) = \alpha, \quad X(0) = X_0,$$

and an admissible set of controls  $B = \{-1, 1\}$  (for Example 9.21), or  $B = [-1, 1]$  (for Example 9.22). The Hamiltonian for this problem is  $H(\lambda, x) = -|\lambda| + x^2/2$ , and for a

smooth approximation of the  $\lambda$ -derivative, e.g.  $H_\lambda^\delta(\lambda, x) = -\tanh(\lambda/\delta)$ , the non-linear system (9.41) becomes

$$\begin{aligned} 0 &= \bar{X}_{n+1} - \bar{X}_n + \Delta t \tanh(\bar{\lambda}_{n+1}/\delta), \\ 0 &= \bar{\lambda}_n - \bar{\lambda}_{n+1} - \Delta t \bar{X}_n. \end{aligned}$$

Newton's method starts with an initial guess  $(\bar{X}_{n+1}^0, \bar{\lambda}_n^0)$ , for all times  $n = 0, \dots, N-1$ , and updates the solution, for some damping factor  $\gamma \in (0, 1]$ , according to

$$\begin{aligned} \bar{X}_{n+1}^{i+1} &= \bar{X}_{n+1}^i - \gamma \Delta \bar{X}_{n+1}^i, \\ \bar{\lambda}_n^{i+1} &= \bar{\lambda}_n^i - \gamma \Delta \bar{\lambda}_n^i, \end{aligned}$$

where the updates comes from solving the sparse Newton system ( $N = 3$  for illustration)

$$\begin{pmatrix} 1 & -1 & & & & & \\ & d_1^i \Delta t & 1 & & & & \\ & & 1 & -\Delta t & -1 & & \\ & & & -1 & d_2^i \Delta t & 1 & \\ & & & & 1 & -\Delta t & \\ & & & & & -1 & 1 \end{pmatrix} \begin{pmatrix} \Delta \bar{\lambda}_0^i \\ \Delta \bar{\lambda}_1^i \\ \Delta \bar{X}_1^i \\ \Delta \bar{\lambda}_2^i \\ \Delta \bar{X}_2^i \\ \Delta \bar{\lambda}_3^i \end{pmatrix} = \begin{pmatrix} \bar{\lambda}_0^i - \bar{\lambda}_1^i - \Delta t \bar{X}_0^i \\ \bar{X}_1^i - \bar{X}_0^i + \Delta t \tanh(\bar{\lambda}_1^i/\delta) \\ \bar{\lambda}_1^i - \bar{\lambda}_2^i - \Delta t \bar{X}_1^i \\ \bar{X}_2^i - \bar{X}_1^i + \Delta t \tanh(\bar{\lambda}_2^i/\delta) \\ \bar{\lambda}_2^i - \bar{\lambda}_3^i - \Delta t \bar{X}_2^i \\ \bar{X}_3^i - \bar{X}_2^i + \Delta t \tanh(\bar{\lambda}_3^i/\delta) \end{pmatrix},$$

and  $d_j^i := \partial_\lambda \tanh(\bar{\lambda}_j^i/\delta) = \delta^{-1} \cosh^{-2}(\bar{\lambda}_j^i/\delta)$ . A Matlab implementation for the above Newton method is shown below, and in Figure 9.9 the solution is shown for different values of  $N$ .

```
% Solving Hamiltonian system with Newton's method
% for T=1, delta=dt and gamma=1

N=1000; dt=1/N;
J=sparse(2*N,2*N); rhs=sparse(2*N,1);
X=sparse(N+1,1); L=sparse(N+1,1);
X(1)= 0.8; % initial data

tol=1;
while tol>1e-6
    % Assemble Newton system row-wise
    for n=1:N
        rhs(2*n-1)=L(n)-L(n+1)-dt*X(n);
        rhs(2*n)=X(n+1)-X(n)+dt*tanh(L(n+1)/dt);
    end
    J(1,1:2)=[1 -1]; J(2*N,2*N-1:2*N)=[-1 1];
    for n=1:N-1
        J(2*n,2*n-1:2*n+1)=[-1 1/cosh(L(n+1)/dt)^2 1];
        J(2*n+1,2*n:2*n+2)=[1 -dt -1];
    end
end
```



```

J(2,1)=0; J(2*N-1,2*N)=0;
% Solve and update
dXL=J\rhs;
L(1)=L(1)-dXL(1); X(N+1)=X(N+1)-dXL(2*N);
for n=2:N
    X(n)=X(n)-dXL(2*n-1); L(n)=L(n)-dXL(2*n-2);
end
tol = norm(rhs) % Error
end

```

### 9.2.3 Convergence of Euler Pontryagin Approximations

**Theorem 9.27.** *Assume that the Hamiltonian  $H$ , defined in (9.4), is Lipschitz continuous on  $\mathbb{R}^d \times \mathbb{R}^d$  and that (9.29) has a solution  $(\bar{X}, \bar{\lambda})$ , where  $\bar{\lambda}_{n+1}$  has uniformly bounded first variation with respect to  $\bar{X}_n$  for all  $n$  and all  $\Delta t$ , that is there is a constant  $K$  such that*

$$|\partial_{\bar{X}_n} \bar{\lambda}_{n+1}| \leq K. \quad (9.43)$$

Then the optimal solution,  $(\bar{X}, \bar{\lambda})$ , of the Pontryagin method (9.29) satisfies the error estimate

$$\begin{aligned} & \left| \inf_{\alpha \in \mathcal{A}} \left( g(X(T)) + \int_0^T h(X(s), \alpha(s)) ds \right) - \left( g(\bar{X}_N) + \Delta t \sum_{n=0}^{N-1} h^\delta(\bar{X}_n, \bar{\lambda}_{n+1}) \right) \right| \\ &= \mathcal{O}(\Delta t + \delta + \frac{(\Delta t)^2}{\delta}) \\ &= \mathcal{O}(\Delta t), \quad \text{for } \delta = \Delta t. \end{aligned} \quad (9.44)$$

The bound  $\mathcal{O}(\Delta t)$  in (9.44) depends on the dimension  $d$  through the Lipschitz norms of the Hamiltonian  $H$  and the constant  $K$  in (9.43).

The work [SS06] presents a convergence result for the case when backward paths  $\bar{X}(t)$  collide on a  $\mathcal{C}^1$  codimension one surface in  $\mathbb{R}^d \times [0, T]$ . The next subsections give a construction of a regularization  $H^\delta$  and the proof of Theorem 9.27.

#### 9.2.3.1 Construction of a Regularization

A possible regularization of  $H$  is to let  $H^\delta$  be a standard convolution mollification of  $H$

$$H^\delta(\lambda, x) = \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} H(z, y) \omega^\delta(z - \lambda) \omega^\delta(y - x) dz dy, \quad (9.45)$$

with  $\omega^\delta : \mathbb{R}^d \rightarrow \mathbb{R}_+$  a  $\mathcal{C}^2$  function compactly supported in the ball  $\{y \in \mathbb{R}^d : |y| \leq \delta\}$  and with integral one  $\int_{\mathbb{R}^d} \omega^\delta(y) dy = 1$ . This regularization remains concave in  $\lambda$ . Our analysis is not dependent of this specific regularization, but uses that

$$\|H - H^\delta\|_C + \delta \|H^\delta\|_{C^1} + \delta^2 \|H^\delta\|_{C^2} = \mathcal{O}(\delta),$$

and that  $H^\delta$  remains a concave function of  $\lambda$ .

### 9.2.3.2 Convergence without Shocks and Colliding Paths

The proof of the theorem is based on four lemmas. In all of those we suppose that the assumptions of Theorem 9.27 are valid.

**Lemma 9.28.** *The discrete dual function is the gradient of the value function, that is*

$$\bar{u}_x(\bar{X}_n, \bar{t}_n) = \bar{\lambda}_n. \quad (9.46)$$

*Proof.* The relation (9.46) holds for  $t_n = T$ . Use the induction assumption that (9.46) holds true for

$t_N \equiv T, t_{N-1}, \dots, t_{n+1}$ . Then the definitions of  $f^\delta$  and  $h^\delta$  imply

$$\begin{aligned} \frac{\partial \bar{u}}{\partial \bar{X}_n}(\bar{X}_n, t_n) &= \partial_{\bar{X}_n}(\bar{u}(\bar{X}_{n+1}, t_{n+1}) + \Delta t h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n)) \\ &= \partial_{\bar{X}_n} \bar{X}_{n+1} \frac{\partial \bar{u}}{\partial \bar{X}_{n+1}}(\bar{X}_{n+1}, t_{n+1}) + \Delta t \partial_{\bar{X}_n} h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) \\ &= (I + \Delta t \partial_{\bar{X}_n} H_\lambda^\delta(\bar{\lambda}_{n+1}, \bar{X}_n)) \bar{\lambda}_{n+1} + \Delta t \partial_{\bar{X}_n} h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) \\ &= \bar{\lambda}_{n+1} + \Delta t \partial_{\bar{X}_n} (H_\lambda^\delta \bar{\lambda} + h^\delta)(\bar{\lambda}_{n+1}, \bar{X}_n) \\ &\quad - \Delta t H_\lambda^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) \partial_{\bar{X}_n} \bar{\lambda}_{n+1} \\ &= \bar{\lambda}_{n+1} + \Delta t H_x^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) \\ &= \bar{\lambda}_n. \end{aligned}$$

□

Section 9.2.7 shows that (9.46) holds precisely for symplectic methods.

We now extend  $\bar{u}$  to be a function defined for all  $t$ . First extend the solution  $\bar{X}$  to all time as a continuous piecewise linear function

$$\bar{X}(t) = \frac{t_{n+1} - t}{\Delta t} \bar{X}_n + \frac{t - t_n}{\Delta t} \bar{X}_{n+1}, \quad \text{for } t_n \leq t < t_{n+1}, \quad (9.47)$$

so that

$$\bar{X}'(t) = H_\lambda^\delta(\bar{\lambda}_{n+1}, \bar{X}_n). \quad (9.48)$$

The following lemma shows that two different solutions can not collide for suitable small  $\Delta t$ .

**Lemma 9.29.** *There is a positive constant  $c$  such that if  $\Delta t \leq c\delta$  two different solutions  $(\bar{X}^1, \bar{\lambda}^1)$  and  $(\bar{X}^2, \bar{\lambda}^2)$  of (9.29) do not intersect.*

*Proof.* Assume there exist two optimal paths  $(\bar{X}^1, \bar{\lambda}^1)$  and  $(\bar{X}^2, \bar{\lambda}^2)$  that intersect at time  $t$ , where  $\bar{t}_n < t \leq \bar{t}_{n+1}$ , then

$$\bar{X}_n^1 + (t - \bar{t}_n) H_\lambda^\delta(\bar{\lambda}_{n+1}^1, \bar{X}_n^1) = \bar{X}_n^2 + (t - \bar{t}_n) H_\lambda^\delta(\bar{\lambda}_{n+1}^2, \bar{X}_n^2)$$

which can be written

$$\bar{X}_n^1 - \bar{X}_n^2 = (t - \bar{t}_n) (H_\lambda^\delta(\bar{\lambda}_{n+1}^2, \bar{X}_n^2) - H_\lambda^\delta(\bar{\lambda}_{n+1}^1, \bar{X}_n^1)). \quad (9.49)$$

To obtain an estimate of the size of the right hand side in (9.49) integrate along the line

$$\bar{X}(s) = \bar{X}_n^1 + s(\bar{X}_n^2 - \bar{X}_n^1),$$

with  $\bar{\lambda}_{n+1}^i$  a function of  $\bar{X}_n^i$ . The difference in the right hand side of (9.49) is

$$\begin{aligned} H_\lambda^\delta(\bar{\lambda}_{n+1}^2, \bar{X}_n^2) - H_\lambda^\delta(\bar{\lambda}_{n+1}^1, \bar{X}_n^1) &= \int_0^1 \frac{dH_\lambda^\delta}{ds} ds \\ &= \int_0^1 (H_{\lambda x}^\delta + H_{\lambda\lambda}^\delta \partial_{\bar{X}_n} \bar{\lambda}_{n+1}) ds (\bar{X}_n^2 - \bar{X}_n^1). \end{aligned}$$

By assumption it holds that  $\|H_{\lambda x}^\delta + H_{\lambda\lambda}^\delta \partial_{\bar{X}_n} \bar{\lambda}_{n+1}\|_C = \mathcal{O}(C_\lambda(1+K)/\delta)$ . Hence the norm of the right hand side in (9.49) is  $\mathcal{O}(\delta^{-1}\Delta t)\mathcal{O}\|\bar{X}_n^1 - \bar{X}_n^2\|$ . Therefore there is a positive constant  $c$  such that if  $\Delta t < c\delta$ , the equation (9.49) has only the solution  $\bar{X}_n^1 = \bar{X}_n^2$ .  $\square$

Since the optimal paths  $\bar{X}$  do not collide, for suitable small  $\Delta t$ , the value function  $\bar{u}$  is uniquely defined along the optimal paths, by (9.31) and

$$\bar{u}(\bar{X}(t), t) = \bar{u}(\bar{X}_{n+1}, t_{n+1}) + (t_{n+1} - t)h^\delta(\bar{X}_n, \bar{\lambda}_{n+1}), \quad t_n < t < t_{n+1} \quad (9.50)$$

and we are ready for the main lemma

**Lemma 9.30.** *The value function for the Pontryagin method satisfies a Hamilton-Jacobi equation close to (9.4), more precisely there holds*

$$\begin{aligned} \bar{u}_t + H(\bar{u}_x, \cdot) &= \mathcal{O}\left(\delta + \Delta t + \frac{(\Delta t)^2}{\delta}\right) \quad \text{in } \mathbb{R}^d \times (0, T), \\ \bar{u} &= g \quad \text{on } \mathbb{R}^d. \end{aligned} \quad (9.51)$$

The error term  $\mathcal{O}(\delta + \Delta t + \frac{(\Delta t)^2}{\delta})$  in (9.51) is a Lipschitz continuous function of  $\bar{u}_x(x, t)$ ,  $x$  and  $t$  satisfying

$$\left|\mathcal{O}\left(\delta + \Delta t + \frac{(\Delta t)^2}{\delta}\right)\right| \leq CC_\lambda \left(\delta + C_x \Delta t + C_x C_\lambda (1+K) \frac{(\Delta t)^2}{\delta}\right),$$

where  $C_x$  and  $C_\lambda$  are the Lipschitz constants of  $H$  in the  $x$  and  $\lambda$  variable, respectively, and  $C \sim 1$  does not depend on the data.

*Proof.* The proof starts with the observation

$$\begin{aligned} 0 &= \frac{d}{dt} \bar{u}(\bar{X}(t), t) + h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) \\ &= \bar{u}_t(\bar{X}(t), t) + \bar{u}_x(\bar{X}(t), t) \cdot f^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) + h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n). \end{aligned} \quad (9.52)$$

The idea is now to use that the dual function  $\bar{\lambda}$  is the gradient of  $\bar{u}$  at the time levels  $t_n$ , by Lemma 9.28, (and a good approximation at times in between) and that the modified discrete Pontryagin method shows that the right hand side in (9.52) is consistent with the correct Hamiltonian  $H$ .

We will first derive an estimate of  $|\bar{u}_x(\bar{X}(t), t) - \bar{\lambda}_{n+1}|$  for  $t_n < t < t_{n+1}$ . We have that

$$\bar{u}(\bar{X}(t), t) = \bar{u}(\bar{X}_{n+1}, \bar{t}_{n+1}) + (\bar{t}_{n+1} - t)h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n)$$

Therefore  $\bar{u}_x(\bar{X}(t), t)$  can be written as

$$\begin{aligned}\bar{u}_x(\bar{X}(t), t) &= \frac{\partial \bar{X}_n}{\partial \bar{X}_t} \left( \frac{\partial \bar{X}_{n+1}}{\partial \bar{X}_n} \bar{u}_x(\bar{X}_{n+1}, t_{n+1}) + (t_{n+1} - t) \partial_{\bar{X}_n} h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) \right) \\ &= \frac{\partial \bar{X}_n}{\partial \bar{X}_t} \left( \frac{\partial \bar{X}_{n+1}}{\partial \bar{X}_n} \bar{\lambda}_{n+1} + (t_{n+1} - t) \partial_{\bar{X}_n} h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) \right).\end{aligned}$$

Introduce the notation

$$\begin{aligned}A &\equiv \partial_{\bar{X}_n} H_\lambda^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) = H_{\lambda x}^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) + H_{\lambda \lambda}^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) \partial_{\bar{X}_n} \bar{\lambda}_{n+1} \\ &= \mathcal{O}(C_\lambda(1 + K)/\delta).\end{aligned}\tag{9.53}$$

We have

$$\begin{aligned}\frac{\partial \bar{X}_{n+1}}{\partial \bar{X}_n} &= I + \Delta t A = I + (t - t_n)A + (t_{n+1} - t)A \\ \frac{\partial \bar{X}_n}{\partial \bar{X}_t} &= (I + (t - t_n)A)^{-1}\end{aligned}$$

therefore as in Lemma 9.28

$$\begin{aligned}\bar{u}_x(\bar{X}(t), t) &= \bar{\lambda}_{n+1} + (t_{n+1} - t)(I + (t - t_n)A)^{-1}(A\bar{\lambda}_{n+1} + \partial_{\bar{X}_n} h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n)) \\ &= \bar{\lambda}_{n+1} + (t_{n+1} - t)(I + (t - t_n)A)^{-1}H_x^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) \\ &= \bar{\lambda}_{n+1} + \mathcal{O}(C_x \Delta t + C_x C_\lambda (K + 1)(\Delta t)^2/\delta).\end{aligned}\tag{9.54}$$

Introduce the notation  $\tilde{\lambda} \equiv \bar{u}_x(\bar{X}(t), t)$  and split the Hamiltonian term in (9.52) into three error parts:

$$\begin{aligned}r(\tilde{\lambda}, \bar{X}(t), t) &\equiv \tilde{\lambda} f^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) + h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) - H(\tilde{\lambda}, \bar{X}(t)) \\ &= \tilde{\lambda} f^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) + h^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) - H^\delta(\tilde{\lambda}, \bar{X}_n) \\ &\quad + H^\delta(\tilde{\lambda}, \bar{X}_n) - H^\delta(\tilde{\lambda}, \bar{X}(t)) \\ &\quad + H^\delta(\tilde{\lambda}, \bar{X}(t)) - H(\tilde{\lambda}, \bar{X}(t)) \\ &\equiv I + II + III.\end{aligned}\tag{9.55}$$

Taylor expansion of  $H^\delta$  to second order and (9.54) show

$$\begin{aligned}|I| &= |H^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) + (\tilde{\lambda} - \bar{\lambda}_{n+1})H_\lambda^\delta(\bar{\lambda}_{n+1}, \bar{X}_n) - H^\delta(\tilde{\lambda}, \bar{X}_n)| \\ &\leq \min(2C_\lambda |\tilde{\lambda} - \bar{\lambda}_{n+1}|, |(\tilde{\lambda} - \bar{\lambda}_{n+1})H_{\lambda \lambda}^\delta(\xi, \bar{X}_n)(\tilde{\lambda} - \bar{\lambda}_{n+1})|/2) \\ &\leq CC_\lambda (C_x \Delta t + C_x C_\lambda (K + 1)(\Delta t)^2/\delta);\end{aligned}$$

the Lipschitz continuity of  $H^\delta$  implies

$$|II| \leq |H_x^\delta| |\bar{X}(t) - \bar{X}_n| \leq |H_x^\delta| |H_\lambda^\delta| \Delta t;$$

and the approximation  $H^\delta$  satisfies

$$|III| \leq CC_\lambda \delta.$$

The combination of these three estimates proves (9.51).

To finish the proof of the lemma we show that the error function  $r$  can be extended to a Lipschitz function in  $\mathbb{R}^d \times \mathbb{R}^d \times [0, T]$ . We note that by (9.43), (9.47) and (9.54)  $\tilde{\lambda}$  is a Lipschitz function of  $X_t$  and  $t$ , and  $r(\tilde{\lambda}(X_t, t), X_t, t)$  is Lipschitz in  $X_t$  and  $t$ . By

$$r(\lambda, X, t) \equiv r(\tilde{\lambda}(X, t), X, t)$$

we obtain a Lipschitz function  $r$  in  $\mathbb{R}^d \times \mathbb{R}^d \times [0, T]$ . □

The results in these lemmas finishes the proof of Theorem 9.27: the combination of the residual estimates in Lemma 9.30 and the  $\mathcal{C}$ -stability estimate of viscosity solutions in Lemma 9.31 proves the theorem.

The approximation result can be extended to the case when the set of backward optimal paths  $\{(\bar{X}(t), t) : t < T\}$ , solving (9.29) and (9.47), may collide into a codimension one surface  $\Gamma$  in space-time  $\mathbb{R}^d \times [0, T]$ , see [SS06].

### 9.2.3.3 Maximum Norm Stability for Hamilton-Jacobi Equations

The seminal construction of viscosity solutions by Crandall and Lions [?] also includes  $\mathcal{C}$  stability results formulated in a general setting. We restate a variant adapted to the convergence results in this paper.

**Lemma 9.31.** *Suppose  $H : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a Lipschitz continuous Hamiltonian satisfying for a constant  $C$  and for all  $x, \hat{x}, \lambda, \hat{\lambda} \in \mathbb{R}^d$*

$$\begin{aligned} |H(\lambda, x) - H(\lambda, \hat{x})| &\leq C_x |x - \hat{x}| (1 + |\lambda|), \\ |H(\lambda, x) - H(\hat{\lambda}, x)| &\leq C_\lambda |\lambda - \hat{\lambda}|. \end{aligned}$$

*Suppose also that  $e : \mathbb{R}^d \times [0, T] \rightarrow \mathbb{R}$  and  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  are Lipschitz continuous. Then, the bounded uniformly continuous viscosity solutions  $u$  and  $\hat{u}$  of the Hamilton-Jacobi equations*

$$u_t + H(u_x, \cdot) = 0 \quad \text{in } \mathbb{R}^d \times (0, T), \quad u|_{\mathbb{R}^d \times \{T\}} = g, \quad (9.56)$$

$$\hat{u}_t + H(\hat{u}_x, \cdot) = e \quad \text{in } \mathbb{R}^d \times (0, T), \quad \hat{u}|_{\mathbb{R}^d \times \{T\}} = g, \quad (9.57)$$

*satisfy the  $\mathcal{C}$ -stability estimate*

$$\mathcal{O} \|u - \hat{u}\|_{\mathcal{C}(\mathbb{R}^d \times [0, T])} \leq T \mathcal{O} \|e\|_{\mathcal{C}(\mathbb{R}^d \times [0, T])}. \quad (9.58)$$

This follows from the maximum norm stability (9.22), but other proofs based on the maximum principle or the comparison principle are also possible, see [SS06].

### 9.2.4 How to obtain the Controls

The optimal control for the exact problem (9.4) is determined by the value function through the Pontryagin principle

$$\alpha(x, t) \in \operatorname{argmin}_{a \in B} (u_x(x, t) \cdot f(x, a) + h(x, a)).$$

Assume we have solved a discrete approximating optimal control problem and obtained the approximations  $\bar{X}$ ,  $\bar{\lambda}$  and  $\bar{u}$ . Can they be used to determine an approximation of the control  $\alpha$ ? Even in the case that the optimal control  $S(\lambda, x) \equiv \operatorname{argmin}_a (\lambda \cdot f(x, a) + h(x, a))$  is a function, it is in general not continuous as function of  $x$  and  $\lambda$  but only piecewise Lipschitz continuous. Therefore the approximate control  $S(\bar{\lambda}(t), x)$  cannot be accurate in maximum norm. However, weaker measures of the control can converge; for instance the value function is accurately approximated in Theorems 9.27 and ???. At the points where  $S$  is Lipschitz continuous the error in the control is proportional to the error  $|\bar{\lambda}(x, t) - u_x(x, t)|$ , for fixed  $x$ . If we assume that the error  $\bar{u}(\cdot, t) - u(\cdot, t)$  is bounded by  $\epsilon$  in a  $\sqrt{\epsilon}$ -neighborhood of  $x$  and that  $\bar{u}_{xx}$  and  $u_{xx}$  also are bounded there, we obtain, for difference quotients  $\Delta u / \Delta x$  and  $|\Delta x| = \sqrt{\epsilon}$ , the error estimate

$$\bar{\lambda} - u_x = \bar{\lambda} - \frac{\Delta \bar{u}}{\Delta x} + \frac{\Delta \bar{u}}{\Delta x} - \frac{\Delta u}{\Delta x} + \frac{\Delta u}{\Delta x} - u_x = \mathcal{O}(\Delta x + \epsilon / \Delta x) = \mathcal{O}(\sqrt{\epsilon}).$$

Convergence of the approximate path  $(\bar{X}, \bar{\lambda})$  typically requires Lipschitz continuous flux  $(H_\lambda, H_x)$ , which we do not assume in this work.

### 9.2.5 Inverse Problems and Tikhonov Regularization

One way to introduce regularization of ill-posed inverse problems is to study a simple example such as  $u' = f$ : the forward problem to determine  $u$  from  $f$  in this case becomes a well-posed integral  $u(x) = u(0) + \int_0^x f(s) ds$  and the inverse problem is then to determine  $f$  from  $u$  by the derivative  $f = u'$ . Note that a small error in the data can be amplified when differentiated; for instance a small perturbation maximum-norm  $\epsilon \sin(\omega x)$  in  $u$  leads to the  $f$ -perturbation  $\epsilon \omega \cos(\omega x)$  which is large (in maximum-norm) if  $\omega \epsilon \gg 1$  even if  $\epsilon \ll 1$ , while a small maximum-norm perturbation of  $f$  leads to a small perturbation of  $u$  (in maximum norm). This is the reason that, to determine  $u$  from  $f$  is well posed (in maximum norm), while the inverse problem to determine  $f$  from  $u$  is ill posed.

The simplest method to regularize the problem  $f = u'$  is to replace the derivative with a difference quotient with suitable step size  $h$ . If we assume that our measured values  $u_*$  of  $u \in \mathcal{C}^2$  are polluted with an error  $\eta$  of size  $\epsilon$  in maximum norm so that  $u_* = u + \eta$ , we have

$$f = (u_* - \eta)'$$

To avoid differentiating  $\eta$  we use the difference quotient

$$\begin{aligned} f(x) &= u'(x) \\ &= \frac{u(x+h) - u(x)}{h} + \mathcal{O}(h) \\ &= \frac{u_*(x+h) - u_*(x)}{h} + \mathcal{O}(\epsilon h^{-1} + h). \end{aligned}$$

The error term is minimal if we choose  $h^2 \simeq \epsilon$ , that is the optimal step size,  $h \simeq \sqrt{\epsilon}$ , yields the error  $\mathcal{O}(\epsilon^{1/2})$  to compute  $f$  by the difference quotient. This difference quotient converges to  $u'$  as  $\epsilon$  tends to zero. If we take too small step size (e.g.  $h = \epsilon$ ), the estimation error does not tend to zero as the measurement error tends to zero.

We can write the inverse problem  $u' = f$  as the optimal control problem

$$\begin{aligned} \dot{X}^t &= \alpha^t, \\ \min_{\alpha: (0,1) \rightarrow [-M, M]} & 2^{-1} \int_0^1 |X^t - X_*^t|^2 dt, \end{aligned}$$

where we changed notation to  $t := x$ ,  $X = u$ ,  $X_* = u_*$ ,  $\alpha := f$  and put the constraint to seek  $\alpha$  in the bounded set  $[-M, M]$  for some positive  $M$ . The Hamiltonian becomes

$$H(\lambda, x, t) = \min_{\alpha \in [-M, M]} (\lambda \cdot \alpha + 2^{-1} |x - X_*^t|^2) = -M|\lambda| + 2^{-1} |x - X_*^t|^2$$

which is not differentiable and leads to the system

$$\begin{aligned} \dot{X}^t &= -M \operatorname{sgn}(\lambda) \\ \dot{\lambda}^t &= -(X^t - X_*^t). \end{aligned}$$

A regularization of this is to replace  $\operatorname{sgn} \lambda$  by  $\tanh \lambda/\delta$  in the flux, which yields the regularized Hamiltonian

$$H^\delta(\lambda, x, t) = -M\delta \log(\cosh \frac{\lambda}{\delta}) + 2^{-1} |x - X_*^t|^2. \quad (9.59)$$

A standard alternative and related regularization is to add a penalty function depending on the control to the Lagrangian

$$\mathcal{L}^\delta(\lambda, x, \alpha) := \int_0^1 \lambda^t (\alpha^t - \dot{X}^t) + 2^{-1} |X^t - X_*^t|^2 + \delta \alpha^2 dt$$

for some  $\delta > 0$ , which generates the Hamiltonian system

$$\begin{aligned} \dot{X}^t &= -M \operatorname{sgn}^\delta(\lambda) \\ \dot{\lambda}^t &= -(X^t - X_*^t), \end{aligned}$$

where  $\operatorname{sgn}^\delta$  is the piecewise linear approximation to  $\operatorname{sgn}$  with slope  $-1/(2\delta)$ , see Figure 9.13. The corresponding Hamiltonian is  $\mathcal{C}^1$  and has the following parabolic approximation

of  $-M|\lambda|$

$$\begin{cases} -\lambda M + \delta M^2 & \text{if } \lambda > 2\delta M \\ -\frac{\lambda^2}{4\delta} & \text{if } -2\delta M \leq \lambda \leq 2\delta M \\ \lambda M + \delta M^2 & \text{if } \lambda \leq -2\delta M, \end{cases}$$

which in some sense is the simplest regularization giving a differentiable Hamiltonian. Such a regularization obtained by adding a penalty function, depending on the control, to the Lagrangian is called a *Tikhonov regularization*. Any smooth modification of the Hamiltonian can be interpreted as adding such a Tikhonov penalty function, see Section 9.2.5. The fundamental property we desire of a regularization is that the Hamiltonian becomes differentiable. It is somewhat difficult to directly see how to choose a penalty yielding differentiable Hamiltonian, therefore we propose instead to directly regularize the Hamiltonian, e.g. by a mollification as in (9.45) (instead of finding appropriate penalty functions):

- choose a suitable set of controls and its range,
- determine the Hamiltonian,
- mollify the Hamiltonian with a parameter  $\delta > 0$  as in (9.45).

Another example of a forward problem is to determine the solution  $u$ , representing e.g. temperature, in the boundary value problem

$$\begin{aligned} (a(x)u'(x))' &= f(x) & 0 < x < 1 \\ u(0) &= u'(1) = 0 \end{aligned} \tag{9.60}$$

for a given source function  $f : (0, 1) \rightarrow (c, \infty)$  and a given conductivity  $a : (0, 1) \rightarrow (c, \infty)$  with  $c > 0$ . This is a well posed problem with the solution

$$u(x) = \int_0^x \frac{F(s) - F(1)}{a(s)} ds,$$

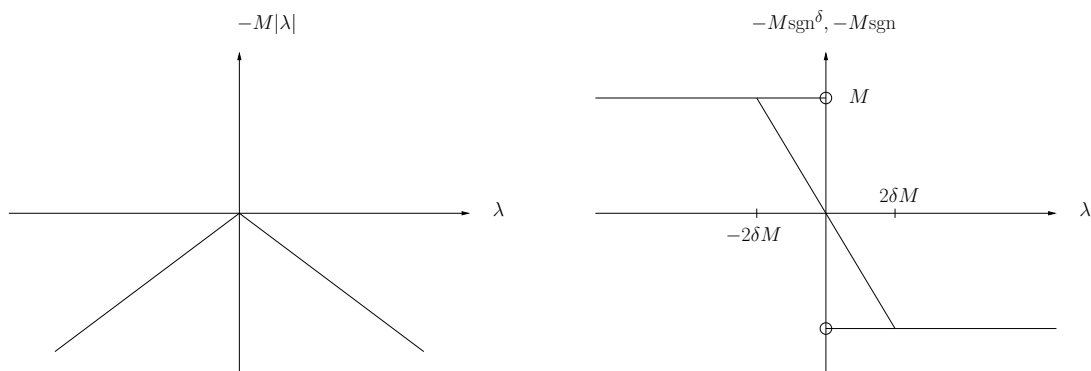


Figure 9.13: Graph of the functions  $-M|\lambda|$ ,  $-\text{sgn}^\delta$  and  $-\text{sgn}$ .



where  $F(s) = \int_0^s f(t) dt$  is a primitive function of  $f$ . The inverse problem to find the conductivity  $a$  from given temperature  $u$  and source  $f$  leads to

$$a(x) = \frac{F(x) - F(1)}{u'(x)}, \quad (9.61)$$

which depends on the derivative  $u'$ , as in the previous example, so that it is ill posed (in maximum norm) by the same reason.

**Example 9.32** (Numerical regularization). Instead of the exact inversion formula (9.61) we can formulate the optimal control problem

$$\min_{a: [0,1] \rightarrow \mathbb{R}} \frac{1}{2} \int_0^1 (u - u^*)^2 + \delta a^2 dx,$$

where  $a$  and  $x$  satisfies (9.60),  $\delta > 0$ , and  $u^*$  denotes given data corresponding to a diffusion coefficient  $a^*$ . From the Lagrangian

$$\begin{aligned} \mathcal{L}(u, \lambda, a) &:= \frac{1}{2} \int_0^1 (u - u^*)^2 + \delta a^2 + (au')' \lambda - f \lambda dx = \\ &= \frac{1}{2} \int_0^1 (u - u^*)^2 + \delta a^2 - au' \lambda' - f \lambda dx, \end{aligned}$$

the Lagrange principle gives that a necessary condition for an optimum is that  $u$ ,  $\lambda$  and  $a$  satisfies Equation (9.60), the dual equation

$$(a(x)\lambda')' = u^* - u, \quad 0 < x < 1, \lambda(0) = \lambda'(1) = 0, \quad (9.62)$$

and

$$u' \lambda' + \delta a = 0, \quad 0 < x < 1. \quad (9.63)$$

In this case the Lagrange principle gives the same result as the Pontryagin principle since the Lagrangian is convex in  $a$ , and since it is smooth in  $a$  no regularization is needed. For  $\delta = 0$ , the Pontryagin principle does not give an explicit Hamiltonian unless we impose some bounds on  $a$ , while the Lagrange principle still is useful numerically, as we shall see.

The simplest way to solve system (9.60), (9.62) and (9.63) is to use the gradient method: given a starting guess  $a^i$ , solve (9.60) to get  $u$ , and (9.62) to get  $\lambda$ , and finally update  $a$  by taking a step of length  $\theta$  in the negative gradient direction, i.e.

$$\begin{aligned} a^{i+1} &= a^i - \theta \frac{d\mathcal{L}(u(a^i), \lambda(a^i), a^i)}{da^i} = a^i - \theta \left( \frac{\partial \mathcal{L}}{\partial u} \frac{du}{da^i} + \frac{\partial \mathcal{L}}{\partial \lambda} \frac{d\lambda}{da^i} + \frac{\partial \mathcal{L}}{\partial a^i} \right) = \\ &= \left\{ \frac{\partial \mathcal{L}}{\partial u} = \frac{\partial \mathcal{L}}{\partial \lambda} = 0 \right\} = a^i - \theta (u' \lambda' + \delta a^i), \quad 0 < x < 1. \end{aligned}$$

Consider the test problem where the measurement  $u^*$  is generated by solving (9.60) with the finite element method for a reference coefficient  $a^*(x) := 1 + 0.5 \sin(2\pi x)$  and a source term  $f = 1$ . To the measurements we add some noise, see Figure 9.14.

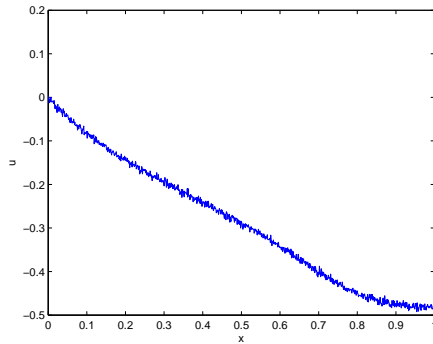


Figure 9.14: Measurements with added noise.

We will now compare different types of regularization: Tikhonov regularization and regularization by discretization or by iteration. In Figure 9.15 the exact inversion (9.61) is shown. A zero misfit error  $u - u^*$  here gives an highly oscillating inversion and is thus infeasible for practical use. The only way to use this method is to introduce a numerical regularization from choosing a sufficiently large discretization. In the right part of Figure 9.15 a 100 times coarser mesh is used for the inversion. It is here possible to see something that vaguely resembles the sought coefficient  $a^*$ .

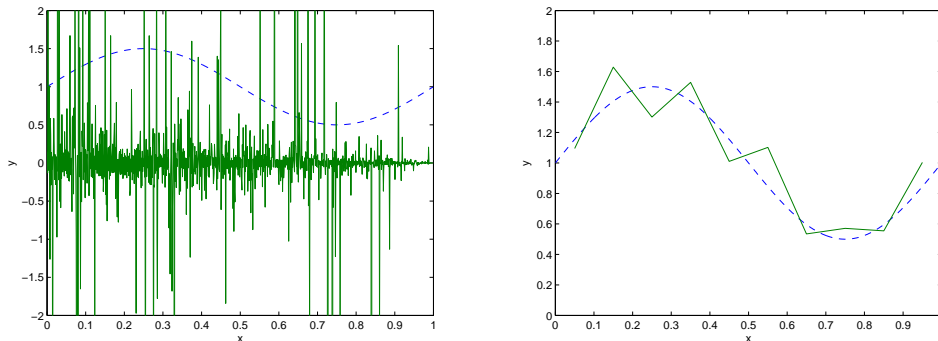


Figure 9.15: Reconstructed coefficient from exact inversion using different meshes.

From the gradient method, for which we choose  $\theta = 10$ , we can in Figure 9.16 see the result for the case with no noise and  $\delta = 0$ . Although the absence of noise will theoretically give an exact fit to data, the method will take a long time to converge, and even for a fast method like Newton's method, a small misfit error may still imply a substantial error in the coefficient.

To test the gradient method for the case with measurement noise we start by letting  $\delta = 0$ . In Figure 9.17 we can see that the gradient method initially finds a smooth function that fits  $\sigma^*$  quite good, but eventually the noise will give a randomly oscillating coefficient as the misfit error decreases. To interrupt the iteration process prematurely

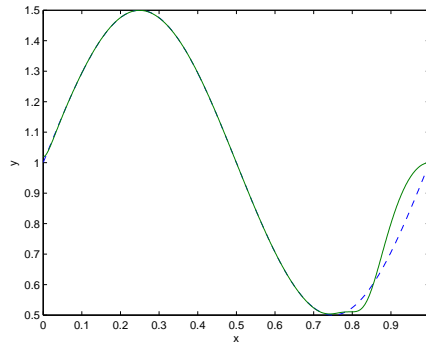


Figure 9.16: Reconstructed coefficient from the gradient method with no noise in measurements and  $\delta = 0$ .

is here a sort of regularization called Landweber iteration [Vog02]. In Figure 9.18 the error in data and coefficients is shown; it is evident that the optimal stopping criterion occurs when the  $\|\sigma - \sigma^*\|$  reaches its minimum. Unfortunately, since  $\sigma^*$  is unknown this criterion cannot be fulfilled in practice.

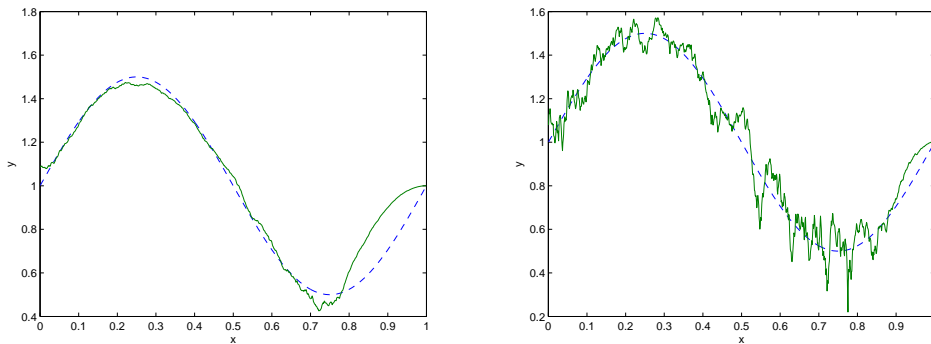


Figure 9.17: Reconstructed coefficient from the gradient method with noisy measurements and  $\delta = 0$ . Left: 100 iterations. Right: 1000 iterations.

In Figure 9.19 the result for the gradient method with a small regularization  $\delta = 5 \cdot 10^{-4}$  is shown. Although the error in the coefficient is higher than for the case with  $\delta = 0$ , in Figure 9.18, this error is bounded and we can thus continue the iterations until the desired tolerance of the gradient norm is met.

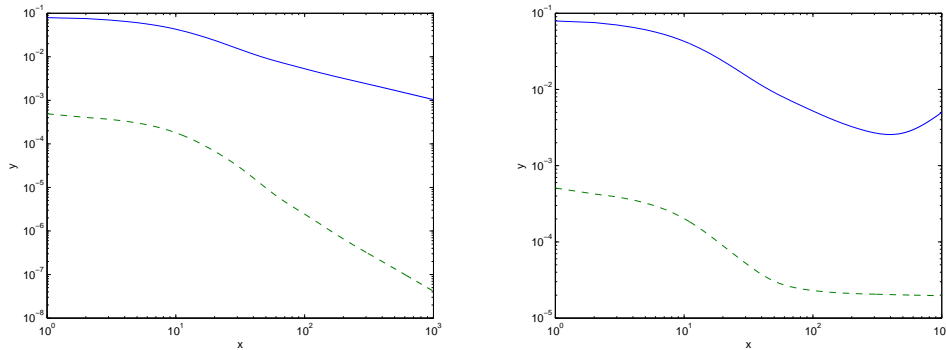


Figure 9.18: Iteration errors from the gradient method. The solid lines depict  $\|\sigma - \sigma^*\|^2$  and the dashed lines show  $\|u - u^*\|^2$ . Left: No noise in data. Right: Noisy data. Note that after a certain number of iterations,  $\|\sigma - \sigma^*\|^2$  will get larger as  $\|u - u^*\|^2$  gets smaller.

**Exercise 9.33.** Consider the the following inverse problems:

- (i) Estimate  $a$  given the solution  $u$  to

$$\begin{aligned} (a(x)u'(x))' &= 1 \quad 0 < x < 1 \\ u(0) &= u(1) = 0. \end{aligned}$$

- (ii) Estimate  $a$  given the boundary solution  $u(1)$  to

$$\begin{aligned} (a(x)u'(x))' &= 0 \quad 0 < x < 1 \\ u(0) &= 0, \\ u'(1) &= 1. \end{aligned}$$

What can we say about the estimation of  $a$  for each problem?

**Example 9.34.** Condition number, matrices, tomography

### 9.2.6 Smoothed Hamiltonian as a Tikhonov Regularization

The  $\mathcal{C}^2$  regularization of the Hamiltonian can also be viewed as a special Tikhonov regularization, using the Legendre transformation: a preliminary idea is to find the Tikhonov penalty function  $T(x, \alpha) : \mathbb{R}^d \times A \rightarrow \mathbb{R}$  such that

$$\min_{\alpha \in A} (\lambda \cdot f(x, \alpha) + T(x, \alpha)) = H^\delta(\lambda, x).$$

In general this can only hold if the set  $A$  is dense enough, e.g. if  $A$  would consist of only two elements the function  $H^\delta$  would not be smooth. Therefore we replace  $A$  seeking the

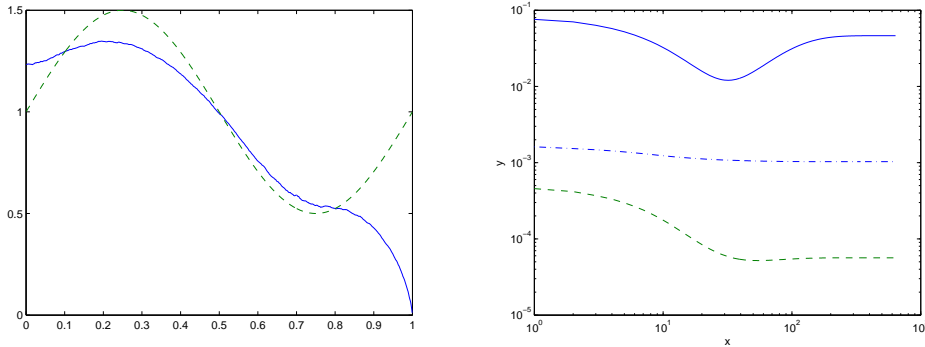


Figure 9.19: Left: Reconstructed coefficient from the gradient method with noisy measurements and  $\delta = 5 \cdot 10^{-4}$ . Right: Errors as in Figure 9.18 but also including the value function  $\|u - u^*\|^2 + \delta \|\sigma - \sigma^*\|^2$  (dash-dotted line).

minimum in the convex closure

$$\widehat{f}(x, A) := \{s f_1 + (1 - s) f_2 \mid s \in [0, 1], \text{ and } f_1, f_2 \in f(x, A)\}$$

and we instead want to find  $T_x(f) : \mathbb{R}^d \times \widehat{f}(x, A) \rightarrow \mathbb{R}$  such that

$$\min_{\phi \in \widehat{f}(x, A)} (\lambda \cdot \phi + T_x(\phi)) = H^\delta(\lambda, x) \quad \text{for all } \lambda \in \mathbb{R}^d. \quad (9.64)$$

To find the Tikhonov penalty, the first step is to observe that by Theorem ?? there is for each  $\lambda$ , where  $\partial_\lambda H(\cdot, x)$  is defined, an  $\alpha$  such that  $\partial_\lambda H(\lambda, x) = f(x, \alpha)$ ; therefore the regularization  $H^\delta(\lambda, x) = \int_{\mathbb{R}^d} H(\lambda - y) \eta(y) dy$ , as in (??), satisfies  $\overline{\partial_\lambda H^\delta(\mathbb{R}^d, x)} \subset \widehat{f}(x, A)$ , since  $H$  is Lipschitz continuous and hence differentiable almost everywhere.

Define the Legendre transformation

$$\tilde{T}_x(\phi) := \sup_{\lambda \in \mathbb{R}^d} \left( -\lambda \cdot \phi + H^\delta(\lambda, x) \right) \quad \text{for all } \phi \in \mathbb{R}^d. \quad (9.65)$$

Figure 9.20 illustrates the value of the Legendre transform

$$T(\phi) = \sup_{\lambda \in \mathbb{R}} \left( -\lambda \cdot \phi + H(\lambda) \right)$$

of a concave differentiable function  $H : \mathbb{R} \rightarrow \mathbb{R}$ , i.e. find the tangent to the curve

$$\{(\lambda, H(\lambda)) \mid \lambda \in \mathbb{R}\}$$

with the slope  $\phi$ , then its intersection with the  $y$ -axis is  $T(\phi)$ ; in multi dimension,  $d \geq 1$ , find the tangent plane of the graph of  $H$  with normal  $(\phi, -1)$ , then the point  $(0, T(\phi))$  is in the plane. If the range of  $\partial_\lambda H(\cdot, x)$  is only a subset  $S$  of  $\mathbb{R}^d$ , we see that  $T(\phi) = +\infty$  for  $\phi \in \mathbb{R}^d - S$ .

**Theorem 9.35.** By defining  $T_x(\phi) := \tilde{T}_x(\phi)$ , the relation (9.64) holds.

*Proof.* Fix a point  $x \in \mathbb{R}^d$ . The definition (9.65) of the Legendre transform implies that for any  $\phi$  and all  $\lambda \in \mathbb{R}^d$  we have

$$\lambda \cdot \phi + \tilde{T}_x(\phi) \geq H^\delta(\lambda, x). \quad (9.66)$$

It remains to show that for any  $\lambda$  we can have equality here by choosing  $\phi$  precisely.

Since the Hamiltonian  $H^\delta(\cdot, x)$  is concave and differentiable, with  $\partial_\lambda H^\delta(\cdot, x) \in \widehat{f}(x, A)$ , the maximum in the Legendre transform is, for  $\phi$  in the interior of  $\widehat{f}(x, A)$ , attained at a point  $\lambda^*$  (depending on  $\phi$ ) satisfying

$$\tilde{T}_x(\phi) = \sup_{\lambda \in \mathbb{R}^d} (-\lambda \cdot \phi + H^\delta(\lambda, x)) = -\lambda^* \cdot \phi + H^\delta(\lambda^*, x)$$

and  $\phi = \partial_\lambda H^\delta(\lambda^*, x)$ , so that the choice  $\phi = \partial_\lambda H^\delta(\lambda, x)$  gives equality in (9.66). The fact that  $\tilde{T}_x$  is lower semicontinuous shows that

$$\inf_{\phi \in \text{interior } \widehat{f}(x, A)} (\lambda \cdot \phi + \tilde{T}_x(\phi)) = \min_{\phi \in \widehat{f}(x, A)} (\lambda \cdot \phi + \tilde{T}_x(\phi)).$$

□

**Exercise 9.36.** Show that Tikhonov penalty for the regularized Hamiltonian (9.59) in the  $u' = f$  problem is

$$\frac{\delta M^2}{2} \left( \left(1 + \frac{\alpha}{M}\right) \log\left(1 + \frac{\alpha}{M}\right) + \left(1 - \frac{\alpha}{M}\right) \log\left(1 - \frac{\alpha}{M}\right) \right) + \frac{1}{2} |x - X_*^t|^2.$$

## 9.2.7 General Approximations

The essential property of the symplectic Euler method we have used is that  $\bar{u}_x(\bar{X}_n, t_n) = \bar{\lambda}_n$ . This relation holds precisely for symplectic approximations (cf. Remark 9.38):

**Theorem 9.37.** Consider a general one step method

$$\begin{aligned} \bar{X}_{n+1} &= A(\bar{\lambda}_{n+1}, \bar{X}_n) \\ \bar{\lambda}_n &= C(\bar{\lambda}_{n+1}, \bar{X}_n) \end{aligned} \quad (9.67)$$

with

$$\bar{u}(\bar{X}_n, t_n) = g(\bar{X}_n) + \sum_{m=n}^{N-1} B(\bar{\lambda}_{m+1}, \bar{X}_m) \Delta t.$$

Then  $\bar{u}_x(\bar{X}_n, t_n) = \bar{\lambda}_n$ , for all  $n$ , implies that the mapping  $\phi : (\bar{X}_n, \bar{\lambda}_n) \mapsto (\bar{X}_{n+1}, \bar{\lambda}_{n+1})$  is symplectic. If  $\phi$  is symplectic it is possible to choose the function  $B$  so that  $\bar{u}_x(\bar{X}_n, t_n) = \bar{\lambda}_n$ , for all  $n$ .

*Proof.* As in Lemma 9.28 we have

$$\bar{u}_x(\bar{X}_n, t_n) = \frac{dA(\bar{X}_n, \bar{\lambda}_{n+1}(\bar{X}_n))}{d\bar{X}_n} \bar{u}_x(\bar{X}_{n+1}, t_{n+1}) + \frac{dB(\bar{X}_n, \bar{\lambda}_{n+1}(\bar{X}_n))}{d\bar{X}_n}.$$

Therefore the relation

$$\bar{u}_x(\bar{X}_n, t_n) = \bar{\lambda}_n$$

holds if and only if  $\lambda A_\lambda + B_\lambda = 0$  and  $\lambda A_x + B_x = C$ . Let  $S \equiv \lambda A + B$ . Then  $\lambda A_\lambda + B_\lambda = 0$  is equivalent to  $S_\lambda = A$ , but  $S_\lambda = A$  implies  $B = S - \lambda S_\lambda$  so that  $\lambda A_x + B_x = S_x$ . Therefore  $\lambda A_\lambda + B_\lambda = 0$  and  $\lambda A_x + B_x = C$  is equivalent to  $A = S_\lambda$  and  $C = S_x$ .

Let  $S \equiv \bar{\lambda}_{n+1} \cdot \bar{X}_n + \Delta t \tilde{H}(\bar{\lambda}_{n+1}, \bar{X}_n)$ . Then (9.67), with  $A = S_\lambda$  and  $C = S_x$ , becomes

$$\begin{aligned} \bar{X}_{n+1} &= \bar{X}_n + \Delta t \tilde{H}_\lambda(\bar{X}_n, \bar{\lambda}_{n+1}) \\ \bar{\lambda}_n &= \bar{\lambda}_{n+1} + \Delta t \tilde{H}_x(\bar{X}_n, \bar{\lambda}_{n+1}), \end{aligned} \tag{9.68}$$

which by Remark 9.38 is equivalent to symplecticity of the mapping  $(\bar{X}_n, \bar{\lambda}_n) \mapsto (\bar{X}_{n+1}, \bar{\lambda}_{n+1})$ .  $\square$

**Remark 9.38.** A one step method (9.67), interpreted as

$$(\bar{X}_n, \bar{\lambda}_n) \mapsto (\bar{X}_{n+1}, \bar{\lambda}_{n+1}),$$

is called symplectic if there exists a function  $\tilde{H}(\bar{\lambda}_{n+1}, \bar{X}_n)$  such that (9.68) holds, see Theorem 5.1, Lemma 5.2 and (5.5) in Chapter VI of [HLW02], where a thorough study on symplectic methods can be found.

To generalize the error estimate of Theorems 9.27 and ?? to general symplectic one step approximations (9.68), e.g. the second order symplectic Runge-Kutta method

$$\tilde{H} = \frac{1}{2} \left( H(\bar{\lambda}_{n+1}, \bar{X}_n) + H(\bar{\lambda}_{n+1} + \Delta t H_x(\bar{\lambda}_{n+1}, \bar{X}_n), \bar{X}_n + \Delta t H_\lambda(\bar{\lambda}_{n+1}, \bar{X}_n)) \right)$$

requires first an extension of  $\bar{X}_n$  and  $\bar{u}$  to all time, by approximations  $(\bar{f}, \bar{h})$  of  $(f^\delta, h^\delta)$  with

$$\frac{d\bar{X}}{dt} = \bar{f} \quad \text{and} \quad \frac{d\bar{u}}{dt} = -\bar{h},$$

and then an estimate of the residual error  $r$  as in (9.55). In practice we need more regularity of  $H^\delta$  to take advantage of higher order methods. Since we only have Lipschitz bounds of  $H$  the estimate of  $r$  is not smaller than the error  $h^\delta - \bar{h}$ , which is  $\mathcal{O}(\|H^\delta\|_{C^p})(\Delta t)^p = \mathcal{O}((\Delta t)^p/\delta^{p-1})$  for a  $p$ th order accurate method. Consequently the residual error is not smaller than  $\mathcal{O}(\delta + (\Delta t)^p/\delta^{p-1}) = \mathcal{O}(\Delta t)$  for  $\delta \simeq \Delta t$ , so that our error estimate does not improve for higher order schemes, without additional assumptions. On the other hand by extending  $\bar{X}$  as a piecewise linear function, as before, the only change of the analysis in Sections 9.2.3.2 and ?? to other symplectic methods (9.68) is to replace  $H^\delta(\bar{\lambda}_{n+1}, \bar{X}_n)$  by  $\tilde{H}(\bar{\lambda}_{n+1}, \bar{X}_n)$  and since

$$\|H^\delta - \tilde{H}\|_C + \delta \|H^\delta - \tilde{H}\|_{C^1} + \delta^2 \|H^\delta - \tilde{H}\|_{C^2} = \mathcal{O}(\Delta t)$$

the estimate (9.51) holds for all symplectic methods which are at least first order accurate.

Similarly, by considering  $(\bar{X}_{n+1}, \bar{\lambda}_n)$ , instead of  $(\bar{X}_n, \bar{\lambda}_{n+1})$ , as independent variables the scheme

$$\begin{aligned}\bar{X}_n &= A(\bar{X}_{n+1}, \bar{\lambda}_n) \\ \bar{\lambda}_{n+1} &= C(\bar{X}_{n+1}, \bar{\lambda}_n),\end{aligned}$$

is symplectic if and only if

$$\begin{aligned}\bar{X}_n &= \bar{X}_{n+1} - \Delta t \hat{H}_\lambda(\bar{X}_{n+1}, \bar{\lambda}_n) \\ \bar{\lambda}_{n+1} &= \bar{\lambda}_n - \Delta t \hat{H}_x(\bar{X}_{n+1}, \bar{\lambda}_n),\end{aligned}\tag{9.69}$$

and the error analysis of the methods (9.68) applies with

$$\tilde{H}(\bar{X}_n, \bar{\lambda}_{n+1}) = (\bar{X}_{n+1} - \bar{X}_n) \cdot (\bar{\lambda}_{n+1} - \bar{\lambda}_n) + \hat{H}(\bar{X}_{n+1}, \bar{\lambda}_n).$$

An example of a method (9.69) is the Euler method  $\hat{H} = H$ , which is backward Euler for  $\bar{X}$  forwards in time and backward Euler for  $\bar{\lambda}$  backwards in time, in contrast to (9.29) which is forward Euler for  $\bar{X}$  forwards in time and forward Euler for  $\bar{\lambda}$  backwards in time.

### 9.3 Optimal Control of Stochastic Differential Equations

In this section we study optimal control of the solution  $X(t)$  to the stochastic differential equation

$$\begin{cases} dX_i &= a_i(X(s), \alpha(s, X(s)))dt + b_{ij}(X(s), \alpha(s, X(s)))dW_j, \quad t < s < T \\ X(t) &= x \end{cases}\tag{9.70}$$

where  $T$  is a fixed terminal time and  $x \in \mathbb{R}^n$  is a given initial point. Assume that  $a_i, b_{ij} : \mathbb{R}^n \times A \rightarrow \mathbb{R}$  are smooth bounded functions, where  $A$  is a given compact subset of  $\mathbb{R}^m$ . The function  $\alpha : [0, T] \times \mathbb{R}^n \rightarrow A$  is a *control* and let  $\mathcal{A}$  be the set of admissible Markov control functions  $t \rightarrow \alpha(t, X(t))$ . The Markov control functions use the current value  $X(s)$  to affect the dynamics of  $X$  by adjusting the drift and the diffusion coefficients. Let us for these admissible controls  $\alpha \in \mathcal{A}$  define the *cost*

$$C_{t,x}(\alpha) = E\left[\int_t^T h(X(s), \alpha(s))ds + g(X(T))\right]$$

where  $X$  solves the stochastic differential equation (9.70) with control  $\alpha$  and

$$h : \mathbb{R}^n \times A \rightarrow \mathbb{R}, \quad g : \mathbb{R}^n \rightarrow \mathbb{R}$$

are given smooth bounded functions. We call  $h$  the *running cost* and  $g$  the *terminal cost*. Our goal is to find an optimal control  $\alpha^*$  which minimizes the expected cost,  $C_{t,x}(\alpha)$ .

Let us define the value function

$$u(t, x) \equiv \inf_{\alpha \in \mathcal{A}} C_{t,x}(\alpha).\tag{9.71}$$



The plan is to show that  $u$  solves a certain Hamilton-Jacobi equation and that the optimal control can be reconstructed from  $u$ . We first assume for simplicity that the optimal control is attained, i.e

$$u(t, x) = \min_{\alpha \in \mathcal{A}} C_{t,x}(\alpha) = C_{t,x}(\alpha^*).$$

The generalization of the proofs without this assumption is discussed in Exercise 9.45.

### 9.3.1 An Optimal Portfolio

**Example 9.39.** Assume that the value of a portfolio,  $X(t)$ , consists of risky stocks,  $S(t) = \alpha(t)X(t)$ , and risk less bonds,  $B(t) = (1 - \alpha(t))X(t)$ , where  $\alpha(t) \in [0, 1]$  and

$$dS = aSdt + cSdW, \quad (9.72)$$

$$dB = bBdt, \quad (9.73)$$

with  $0 \leq b < a$ . Define for a given function  $g$  the cost function

$$C_{t,x}(\alpha) = E[g(X(T)) | X(t) = x].$$

Then our goal is to determine the Markov control function  $\alpha(t, X(t))$ , with  $\alpha : [0, T] \times \mathbb{R} \rightarrow [0, 1]$  that maximizes the cost function. The solution will be based on the function

$$u(t, x) \equiv \max_{\alpha} C_{t,x}(\alpha),$$

and we will show that  $u(t, x)$  satisfies the following *Hamilton-Jacobi* equation,

$$u_t + \max_{\alpha \in [0,1]} \left\{ (a\alpha + b(1 - \alpha))xu_x + \frac{c^2\alpha^2}{2}x^2u_{xx} \right\} = 0, \quad (9.74)$$

$$u(T, x) = g(x),$$

that is

$$u_t + H(x, u_x, u_{xx}) = 0$$

for

$$H(x, p, w) \equiv \max_{v \in [0,1]} (av + b(1 - v)xp + \frac{c^2v^2}{2}x^2w).$$

**Example 9.40.** Assume that  $u_{xx} < 0$  in the equation (9.74). Determine the optimal control function  $\alpha^*$ .

**Solution.** By differentiating  $f(\alpha) = (a\alpha + b(1 - \alpha))xu_x + \frac{c^2\alpha^2}{2}x^2u_{xx}$  in (9.74) with respect to  $\alpha$  and using  $df/d\alpha = 0$ , we obtain

$$\hat{\alpha} = -\frac{(a - b)u_x}{c^2xu_{xx}}.$$

Then the optimal control  $\alpha^*$  is given by

$$\alpha^* = \begin{cases} 0, & \text{if } \hat{\alpha} < 0 \\ \hat{\alpha}, & \text{if } \hat{\alpha} \in [0, 1] \\ 1 & \text{if } 1 < \hat{\alpha} \end{cases}$$

The optimal value yields in (9.74) the Hamilton-Jacobi equation

$$u_t + H(x, u_x, u_{xx}) = 0,$$

where

$$H(x, u_x, u_{xx}) = \begin{cases} bxu_x, & \text{if } \hat{\alpha} < 0 \\ bxu_x - \frac{(a-b)^2 u_x^2}{2c^2 u_{xx}}, & \text{if } \hat{\alpha} \in [0, 1] \\ axu_x + \frac{c^2 x^2 u_{xx}}{2} & \text{if } 1 < \hat{\alpha} \end{cases} \quad (9.75)$$

□

**Example 9.41.** What is the optimal control function  $\alpha = \alpha^*$  for  $g(x) = x^r, 0 < r < 1$  ?

**Solution.** We have  $dX = d(\alpha X + (1 - \alpha)X) = dS + dB = (aS + bB)dt + cSdW = (a\alpha X + b(1 - \alpha)X)dt + c\alpha X dW$ , so that the Itô formula yields

$$\begin{aligned} dg(X) &= dX^r = rX^{r-1}dX + \frac{r(r-1)}{2}X^{r-2}(dX)^2 \\ &= rX^r(a\alpha + b(1 - \alpha))dt + rX^r\alpha c dW + \frac{1}{2}\alpha^2 c^2 r(r-1)X^r dt. \end{aligned}$$

Taking the expectation value in the above,

$$E[X^r(T)] = X^r(0) + E \left[ \int_0^T rX^r \left( a\alpha + b(1 - \alpha) + \frac{1}{2}\alpha^2 c^2 (r-1) \right) dt \right].$$

Finally, perturb the above equation with respect to  $\epsilon \in \mathbb{R}_+$  provided  $\alpha = \alpha^* + \epsilon v$  for some feasible function  $v$ , that is  $\alpha^* + \epsilon v \in [0, 1]$  for sufficiently small  $\epsilon$ . Then the optimal control,  $\alpha^*$ , should satisfy  $E[X_{\alpha^* + \epsilon v}^r(T)] - E[X_{\alpha^*}^r(T)] \leq 0 \forall v$ . If we make the assumption  $\alpha^* \in (0, 1)$ , then we obtain

$$E \left[ \int_0^T rX^r v (a - b + \alpha^* c^2 (r-1)) dt \right] = 0, \forall v$$

which implies

$$\alpha^* = \frac{a - b}{c^2(1 - r)}.$$

□

**Exercise 9.42.** What is the optimal control in (9.74) for  $g(x) = \log x$  ?

### 9.3.2 Dynamic Programming and Hamilton-Jacobi Equations

**Lemma 9.43.** *Assume that the assumptions in section 9.3.1 hold. Then, the function  $u$  satisfies, for all  $\delta > 0$ , the dynamic programming relation*

$$u(t, x) = \min_{\alpha: [t, t+\delta] \rightarrow \mathcal{A}} E \left[ \int_t^{t+\delta} h(X(s), \alpha(s), X(s)) ds + u(t + \delta, X(t + \delta)) \right]. \quad (9.76)$$

*Proof.* The proof has two steps: to use the optimal control to verify

$$u(t, x) \geq \min_{\alpha \in \mathcal{A}} E \left[ \int_t^{t+\delta} h(X(s), \alpha(s)) ds + u(t + \delta, X(t + \delta)) \right],$$

and then to show that an arbitrary control yields

$$u(t, x) \leq \min_{\alpha \in \mathcal{A}} E \left[ \int_t^{t+\delta} h(X(s), \alpha(s)) ds + u(t + \delta, X(t + \delta)) \right],$$

which together imply Lemma 9.43.

**Step 1:** Choose the optimal control  $\alpha^*$ , from  $t$  to  $T$ , to obtain

$$\begin{aligned} u(t, x) &= \min_{\alpha \in \mathcal{A}} E \left[ \int_t^T h(X(s), \alpha(s), X(s)) ds + g(X(T)) \right] \\ &= E \left[ \int_t^{t+\delta} h(X(s), \alpha^*(s)) ds \right] + E \left[ \int_{t+\delta}^T h(X(s), \alpha^*(s)) ds + g(X(T)) \right] \\ &= E \left[ \int_t^{t+\delta} h(X(s), \alpha^*(s)) ds \right] \\ &\quad + E \left[ E \left[ \int_{t+\delta}^T h(X(s), \alpha^*(s)) ds + g(X(T)) \mid X(t + \delta) \right] \right] \\ &\geq E \left[ \int_t^{t+\delta} h(X(s), \alpha^*(s)) ds \right] + E[u(X(t + \delta), t + \delta)] \\ &\geq \min_{\alpha \in \mathcal{A}} E \left[ \int_t^{t+\delta} h(X(s), \alpha(s), X(s)) ds + u(X(t + \delta), t + \delta) \right]. \end{aligned}$$

**Step 2:** Choose the control  $\alpha^+$  to be arbitrary from  $t$  to  $t + \delta$  and then, given the value  $X(t + \delta)$ , choose the optimal  $\alpha^*$  from  $t + \delta$  to  $T$ . Denote this control by  $\alpha' = (\alpha^+, \alpha^*)$ .

Definition (9.71) shows

$$\begin{aligned}
u(t, x) &\leq C_{t,x}(\alpha') \\
&= E\left[\int_t^T h(X(s), \alpha'(s))ds + g(X(T))\right] \\
&= E\left[\int_t^{t+\delta} h(X(s), \alpha^+(s))ds\right] + E\left[\int_{t+\delta}^T h(X(s), \alpha^*(s))ds + g(X(T))\right] \\
&= E\left[\int_t^{t+\delta} h(X(s), \alpha^+(s))ds\right] \\
&\quad + E\left[E\left[\int_{t+\delta}^T h(X(s), \alpha^*(s))ds + g(X(T))\right] \mid X(t+\delta)\right] \\
&= E\left[\int_t^{t+\delta} h(X(s), \alpha^+(s))ds\right] + E[u(X(t+\delta), t+\delta)].
\end{aligned}$$

Taking the minimum over all controls  $\alpha^+$  yields

$$u(t, x) \leq \min_{\alpha^+ \in \mathcal{A}} E\left[\int_t^{t+\delta} h(X(s), \alpha^+(s))ds + u(X(t+\delta), t+\delta)\right].$$

□

**Theorem 9.44.** *Assume that  $X$  solves (9.70) with a Markov control function  $\alpha$  and that the function  $u$  defined by (9.71) is bounded and smooth. Then  $u$  satisfies the Hamilton-Jacobi equation*

$$\begin{aligned}
u_t + H(t, x, Du, D^2u) &= 0, \\
u(T, x) &= g(x),
\end{aligned}$$

with the Hamiltonian function

$$H(t, x, Du, D^2u) \equiv \min_{\alpha \in \mathcal{A}} \left[ a_i(x, \alpha) \partial_{x_i} u(t, x) + \frac{b_{ik}(x, \alpha) b_{jk}(x, \alpha)}{2} \partial_{x_i x_j} u(t, x) + h(x, \alpha) \right]$$

*Proof.* The proof has two steps: to show that the optimal control  $\alpha = \alpha^*$  yields

$$u_t + a_i^* \partial_{x_i} u + \frac{b_{ik}^* b_{jk}^*}{2} \partial_{x_i x_j} u + h^* = 0, \quad (9.77)$$

where  $a^*(x) = a(x, \alpha^*(t, x))$ ,  $b^*(x) = b(x, \alpha^*(t, x))$  and  $h^*(t, x) = h(t, x, \alpha^*(t, x))$ , and then that an arbitrary control  $\alpha^+$  implies

$$u_t + a_i^+ \partial_{x_i} u + \frac{b_{ik}^+ b_{jk}^+}{2} \partial_{x_i x_j} u + h^+ \geq 0, \quad (9.78)$$

where  $a^+(x) = a(x, \alpha^+(t, x))$ ,  $b^+(x) = b(x, \alpha^+(t, x))$  and  $h^+(t, x) = h(t, x, \alpha^+(t, x))$ . The two equations (9.77) and (9.78) together imply Theorem 9.44.

**Step 1 :** Choose  $\alpha = \alpha^*$  to be the optimal control in (9.70). Then by the dynamic programming principle of Lemma 9.71

$$u(X(t), t) = E\left[\int_t^{t+\delta} h(X(s), \alpha^*(s, X(s)))ds + u(X(t+\delta), t+\delta)\right],$$

so that Itô's formula implies

$$\begin{aligned} -h(t, x, \alpha^*(t, x))dt &= E[du(X(t), t) | X(t) = x] \\ &= (u_t + a_i^* \partial_{x_i} u + \frac{b_{ik}^* b_{jk}^*}{2} \partial_{x_i x_j} u)(t, x)dt. \end{aligned} \quad (9.79)$$

Definition (9.71) shows

$$u(T, x) = g(x),$$

which together with (9.79) prove (9.77).

**Step 2 :** Choose the control function in (9.70) to be arbitrary from time  $t$  to  $t + \delta$  and denote this choice by  $\alpha = \alpha^+$ . The function  $u$  then satisfies by Lemma 9.71

$$u(t, x) \leq E\left[\int_t^{t+\delta} h(X(s), \alpha^+(s))ds\right] + E[u(X(t+\delta), t+\delta)].$$

Hence  $E[du] \geq -h(x, \alpha^+)dt$ . We know that for any given  $\alpha^+$ , by Itô's formula,

$$E[du(t, X(t))] = E\left[u_t + a_i^+ \partial_{x_i} u + \frac{b_{ik}^+ b_{jk}^+}{2} \partial_{x_i x_j} u\right] dt.$$

Therefore, for any control  $\alpha^+$ ,

$$u_t + a_i^+ \partial_{x_i} u + \frac{b_{ik}^+ b_{jk}^+}{2} \partial_{x_i x_j} u + h(x, \alpha^+) \geq 0,$$

which proves (9.78) □

**Exercise 9.45.** Use a minimizing sequence  $\alpha_i$  of controls, satisfying

$$u(t, x) = \lim_{i \rightarrow \infty} C_{t,x}(\alpha_i),$$

to prove Lemma 9.71 and Theorem 9.44 without the assumption that the minimum control is attained.

**Exercise 9.46.** Let  $\mathcal{A}^+$  be the set of all adapted controls  $\{\alpha : [0, T] \times \mathcal{C}[0, T] \rightarrow A\}$  where  $\alpha(s, X)$  may depend on  $\{X(\tau) : \tau \leq s\}$ . Show that the minimum over all adapted controls in  $\mathcal{A}^+$  is in fact the same as the minimum over all Markov controls, that is

$$\inf_{\alpha \in \mathcal{A}^+} C_{t,x}(\alpha) = \inf_{\alpha \in \mathcal{A}} C_{t,x}(\alpha),$$

e.g. by proving the dynamic programming relation (9.76) for adapted controls and motivate why this is sufficient.

### 9.3.3 Relation of Hamilton-Jacobi Equations and Conservation Laws

In this section we will analyze qualitative behavior of Hamilton-Jacobi equations, in particular we will study the limit corresponding to vanishing noise in control of stochastic differential equations. The study uses the relation between the Hamilton-Jacobi equation for  $V : [0, T] \times \mathbb{R} \rightarrow \mathbb{R}$

$$V_t + H(V_x) = 0, \quad V(0, x) = V_0(x), \quad (H - J)$$

and the conservation law for  $U : [0, T] \times \mathbb{R} \rightarrow \mathbb{R}$

$$U_t + H(U)_x = 0, \quad U(0, x) = U_0(x). \quad (C - L)$$

Observe that the substitution  $V(t, x) = \int_{-\infty}^x U(t, y) dy$ , so that  $U = V_x$ , and integration in  $x$  from  $-\infty$  to  $x$  in (C-L) shows

$$V_t + H(V_x) = H(U(t, -\infty)). \quad (9.80)$$

Combined with the assumptions  $U(t, x) \rightarrow 0$  as  $|x| \rightarrow \infty$  and  $H(0) = 0$  we conclude that  $V$  solves (H-J), if  $U$  solves (C-L).

The next step is to understand the nature of the solutions of (C-L). Consider the special Burger's conservation law

$$0 = U_t + U U_x = U_t + \left(\frac{U^2}{2}\right)_x, \quad U(0, x) = U_0(x). \quad (9.81)$$

Let us define a *characteristic path*  $X : [0, T] \times \mathbb{R} \rightarrow \mathbb{R}$  by

$$\frac{dX}{dt}(t) = U(t, X(t)), \quad X(0) = x_0. \quad (9.82)$$

Thus, if  $\psi(t) \equiv U(t, X(t))$  then  $\frac{d\psi}{dt}(t) = 0$  by virtue of (9.81). This means that the value of  $U$  is constant along a characteristic path. If the characteristics do not collide into each other we may expect to find a solution using the initial data  $U_0(x)$  and the set of characteristics. Unfortunately, this is not what happens in general, and collisions between characteristics do exist and give birth to discontinuities known as shocks. For example, this is the case when  $U_0(x) = -\arctan(x)$  and  $t \geq 1$ .

**Exercise 9.47.** Show that  $w(t) = U_x(X(t), t)$  satisfies  $w(t) = w(0)/(1 + w(0)t)$ ,  $t < 1$ , for Burger's equation (9.81) with initial data  $U(x, 0) = -\arctan(x)$ . Hence,  $w(1) = \infty$ , for  $X(0) = 0$ .

Since the method of characteristics does not work globally we have to find an alternative way to explain what happens with the solution  $U(t, x)$  near a shock. It is not enough with the concept of strong or classical solution, since the solution  $U(t, x)$  is not differentiable in general. For this purpose, we define the notion of weak solution. Let  $V$  be the set of test functions  $\{\varphi : (0, +\infty) \times \mathbb{R} \rightarrow \mathbb{R}\}$  which are differentiable and take the

value zero outside some compact set. Then an integrable function  $U$  is a weak solution of (9.81) if it satisfies

$$\int_0^{+\infty} \int_{-\infty}^{+\infty} \left( U(t, x) \varphi_t(t, x) + \frac{U^2(t, x)}{2} \varphi_x(t, x) \right) dx dt = 0, \quad \forall \varphi \in V \quad (9.83)$$

and

$$\int_{-\infty}^{+\infty} |U(t, x) - U_0(x)| dx \rightarrow 0, \quad \text{as } t \rightarrow 0 \quad (9.84)$$

**Example 9.48.** The shock wave

$$U(t, x) = \begin{cases} 1 & x < \frac{t}{2}, \\ 0 & \text{otherwise.} \end{cases}$$

is a weak solution satisfying (9.83) and (9.84). Observe that for  $s \equiv 1/2$

$$\partial_t \int_a^b U dx = \frac{U^2(t, a) - U^2(t, b)}{2} = - \left[ \frac{U^2}{2} \right],$$

and

$$\partial_t \int_a^b U dx = \partial_t ((s t - a) U_- + (b - s t) U_+) = -s(U_+ - U_-),$$

where

$$[w(x_0)] \equiv w_+(x_0) - w_-(x_0) \equiv \lim_{y \rightarrow 0^+} w(x_0 + y) - w(x_0 - y)$$

is the jump at the point  $x_0$ . Consequently, the speed  $s$  of a shock can be determined by the so called *Rankine Hugoniot* condition

$$s[U] = \left[ \frac{U^2}{2} \right]. \quad (9.85)$$

**Exercise 9.49.** Verify that the shock wave solution

$$U_I(t, x) = \begin{cases} 0 & x > -\frac{t}{2}, \\ -1 & \text{otherwise} \end{cases}$$

and the rarefaction wave solution

$$U_{II}(t, x) = \begin{cases} 0 & x \geq 0, \\ \frac{x}{t} & -t < x < 0, \\ -1 & \text{otherwise} \end{cases}$$

are both weak solutions of  $U_t + U U_x = 0$  with the same initial condition.

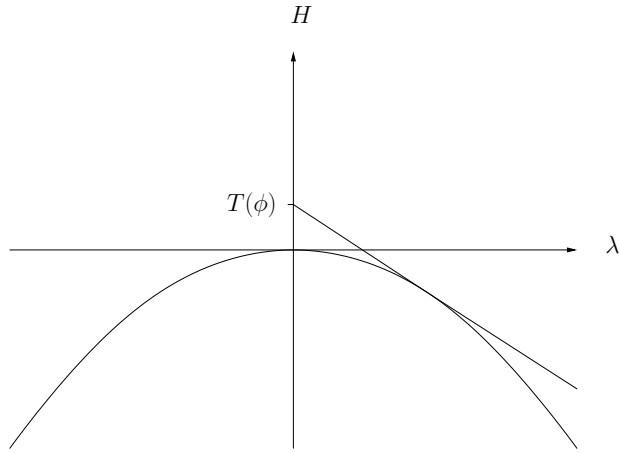


Figure 9.20: Illustration of the Legendre transform. If  $H$  decreases sufficiently fast as  $|\lambda| \rightarrow \infty$ , then  $\partial_\lambda H$  can attain all values in  $\mathbb{R}$  and the range of  $T$  is  $[0, \infty)$ , since  $T(0) = 0$  here. If, on the other hand, the slope of  $H$  is in an interval  $I$ , then  $T(I) = [0, T_+)$  for some upper bound  $T_+$ , and  $T(\mathbb{R} - I) = \{+\infty\}$ .

Figure 9.21: Left: Initial condition. Right: Colliding characteristics and a shock.

Figure 9.22: Shock velocity and Rankine Hugoniot condition

Figure 9.23:  $U_I(t, x)$

Figure 9.24:  $U_{II}(t, x)$



The last exercise shows that we pay a price to work with weak solutions: the lack of uniqueness. Therefore, we need some additional physical information to determine a unique weak solution. This leads us to the concept of *viscosity limit* or *viscosity solution*: briefly, it says that the weak solution  $U$  we seek is the limit  $U = \lim_{\epsilon \rightarrow 0^+} U^\epsilon$  of the solution of the regularized equation

$$U_t^\epsilon + U^\epsilon U_x^\epsilon = \epsilon U_{xx}^\epsilon, \quad \epsilon > 0. \quad (9.86)$$

This regularized equation has continuous and smooth solutions for  $\epsilon > 0$ . With reference to the previous example, the weak solution  $U_{II}$  satisfies  $U_{II} = \lim_{\epsilon \rightarrow 0^+} U^\epsilon$ , but  $U_I \neq \lim_{\epsilon \rightarrow 0^+} U^\epsilon$ . Since a solution of the conservation law can be seen as the derivative of the solution of a Hamilton-Jacobi equation, the same technique of viscosity solutions can be applied to

$$V_t^\epsilon + \frac{(V_x^\epsilon)^2}{2} = \epsilon V_{xx}^\epsilon, \quad \epsilon > 0. \quad (9.87)$$

The functions  $V_I(x, t) = -\int_x^\infty U_I(y, t) dy$ , and  $V_{II}(x, t) = -\int_x^\infty U_{II}(y, t) dy$  have the same initial data and they are both candidates of solutions to the Hamilton-Jacobi equation

$$V_t + \frac{(V_x)^2}{2} = 0.$$

The shock waves for conservation laws corresponds to solutions with discontinuities in the derivative for Hamilton-Jacobi solutions. Only the function  $V_{II}$  satisfies

$$V_{II} = \lim_{\epsilon \rightarrow 0^+} V^\epsilon, \quad (9.88)$$

but  $V_I \neq \lim_{\epsilon \rightarrow 0^+} V^\epsilon$ . It can be shown that the condition (9.88) implies uniqueness for Hamilton-Jacobi equations. Note that (9.88) corresponds to the limit of vanishing noise in control of stochastic differential equations.

### 9.3.4 Numerical Approximations of Conservation Laws and Hamilton-Jacobi Equations

We have seen that the viscous problem

$$\begin{aligned} \partial_t u^\epsilon + \partial_x H(u^\epsilon) &= \epsilon u_{xx}^\epsilon & \text{for } (x, t) \in \mathbb{R} \times (0, +\infty), \\ u^\epsilon(x, 0) &= u_0(x) & \text{for } x \in \mathbb{R}, \end{aligned} \quad (9.89)$$

can be used to construct unique solutions to the conservation law

$$\begin{aligned} \partial_t u + \partial_x H(u) &= 0 & \text{for } (x, t) \in \mathbb{R} \times (0, +\infty), \\ u(x, 0) &= u_0(x) & \text{for } x \in \mathbb{R}. \end{aligned} \quad (9.90)$$

In this section we will develop numerical approximations to the conservation law (9.90) and the related Hamilton-Jacobi equation

$$\partial_t v + H(\partial_x v) = 0,$$

based on viscous approximations. We will also see that too little viscosity may give unstable approximations.

To show the difficulties to solve numerically a problem like (9.90) and (9.89) we consider a related steady-state problem (i.e. a problem that has no dependence on  $t$ )

$$\begin{aligned} \partial_x w(x) - \varepsilon \partial_x^2 w(x) &= 0 \quad \text{for } x < 0, \\ \lim_{x \rightarrow -\infty} w(x) &= 1, \quad w(0) = 0, \end{aligned} \quad (9.91)$$

where  $\varepsilon \geq 0$  is fixed. It is easy to verify that the exact solution is  $w(x) = 1 - \exp(\frac{x}{\varepsilon})$ , for  $x \leq 0$ . Now, we construct a uniform partition of  $(-\infty, 0]$  with nodes  $x_j = j\Delta x$  for  $j = 0, -1, -2, \dots$ , where  $\Delta x > 0$  is a given mesh size. Denoting by  $W_j$  the approximation of  $w(x_j)$ , the use of a second order accurate finite element method or finite difference scheme method leads to the scheme

$$\begin{aligned} \frac{W_{j+1} - W_{j-1}}{2\Delta x} - \varepsilon \frac{W_{j+1} - 2W_j + W_{j-1}}{(\Delta x)^2} &= 0, \quad j = -N + 1, \dots, -1, \\ W_0 &= 0, \\ W_{-N} &= 1. \end{aligned} \quad (9.92)$$

Assume that  $N$  is odd. If  $\varepsilon \ll \Delta x$ , the solution of (9.92) is approximated by

$$\frac{W_{j+1} - W_{j-1}}{2\Delta x} = 0,$$

which yields the oscillatory solution  $W_{2i} = 0$  and  $W_{2i+1} = 1$  that does not approximate  $w$ , instead  $\|w - W\|_{L^2} = \mathcal{O}(1)$ . One way to overcome this difficulty is to replace, in (9.92), the *physical diffusion*  $\varepsilon$  by the *artificial diffusion*  $\hat{\varepsilon} = \max\{\varepsilon, \frac{\Delta x}{2}\}$ . For the general problem  $\beta \cdot \nabla u - \varepsilon \Delta u = f$  take  $\hat{\varepsilon} = \max\{\varepsilon, |\beta| \frac{\Delta x}{2}\}$ . Now, when  $\varepsilon \ll \Delta x$ , we have  $\hat{\varepsilon} = \frac{\Delta x}{2}$  and the method (9.92), with  $\varepsilon$  replaced by  $\hat{\varepsilon}$ , yields  $W_j = W_{j-1}$  for  $j = -N + 1, \dots, -1$ , that is  $W_j = 1$  for  $j = -N, \dots, -1$ , which is an acceptable solution with  $\|w - W\|_{L^2} = \mathcal{O}(\sqrt{\Delta x})$ . Another way to cure the problem is to resolve by choosing  $\Delta x$  small enough, so that  $\hat{\varepsilon} = \varepsilon$ .

The Lax-Friedrich method for the problem (9.90), is given by

$$U_j^{n+1} = U_j^n - \Delta t \left[ \frac{H(U_{j+1}^n) - H(U_{j-1}^n)}{2\Delta x} - \frac{(\Delta x)^2}{2\Delta t} D_+ D_- U_j^n \right], \quad (9.93)$$

with

$$D_+ V_j = \frac{V_{j+1} - V_j}{\Delta x}, \quad D_- V_j = \frac{V_j - V_{j-1}}{\Delta x} \quad \text{and} \quad D_+ D_- V_j = \frac{V_{j+1} - 2V_j + V_{j-1}}{(\Delta x)^2}.$$

The stability condition for the method (9.93) is

$$\lambda \equiv \frac{\Delta x}{\Delta t} > \max_u |H'(u)|. \quad (9.94)$$

We want to approximate the viscosity solution of the one-dimensional Hamilton-Jacobi equation

$$\partial_t v + H(\partial_x v) = 0, \quad (9.95)$$

where  $v = \lim_{\varepsilon \rightarrow 0^+} v^\varepsilon$  and

$$\partial_t v^\varepsilon + H(\partial_x v^\varepsilon) = \varepsilon \partial_x^2 v^\varepsilon. \quad (9.96)$$

Setting  $u = \partial_x v$  and taking derivatives in (9.95), we obtain a conservation law for  $u$ , that is

$$\partial_t u + \partial_x H(u) = 0. \quad (9.97)$$

To solve (9.95) numerically, a basic idea is to apply (9.93) on (9.97) with  $U_i^n = (V_{i+1}^n - V_{i-1}^n)/(2\Delta x)$  and then use summation over  $i$  to approximate the integration in (9.80). We get

$$\begin{aligned} \frac{V_{j+1}^{n+1} - V_{j-1}^{n+1}}{2\Delta x} &= \frac{V_{j+1}^n - V_{j-1}^n}{2\Delta x} \\ &- \Delta t \left[ \frac{H\left(\frac{V_{j+2}^n - V_j^n}{2\Delta x}\right) - H\left(\frac{V_j^n - V_{j-2}^n}{2\Delta x}\right)}{2\Delta x} - \frac{(\Delta x)^2}{2\Delta t} D_+ D_- \frac{V_{j+1}^n - V_{j-1}^n}{2\Delta x} \right]. \end{aligned}$$

Summing over  $j$  and using that  $V_{-\infty}^m = 0$  and  $H(0) = 0$ , it follows that

$$V_j^{n+1} = V_j^n - \Delta t \left[ H\left(\frac{V_{j+1}^n - V_{j-1}^n}{2\Delta x}\right) - \frac{(\Delta x)^2}{2\Delta t} D_+ D_- V_j^n \right], \quad (9.98)$$

which is the Lax-Friedrich method for (9.95). Note that (9.98) is a second order accurate central difference approximation of the equation

$$\partial_t v + H(\partial_x v) = \frac{(\Delta x)^2}{2\Delta t} (1 - (\frac{\Delta t}{\Delta x} H')^2) \partial_x^2 v,$$

which is (9.96) with artificial diffusion  $\Delta x(\lambda^2 - (H')^2)/(2\lambda)$ .

In the two-dimensional case a first order Hamilton-Jacobi equation has the form

$$\partial_t v + H(\partial_{x_1} v, \partial_{x_2} v) = 0. \quad (9.99)$$

The analogous scheme to (9.98) for that equation is

$$\begin{aligned} V_{j,k}^{n+1} = V_{j,k}^n - \Delta t \left[ & H\left(\frac{V_{j+1,k}^n - V_{j-1,k}^n}{2\Delta x_1}, \frac{V_{j,k+1}^n - V_{j,k-1}^n}{2\Delta x_2}\right) \right. \\ & - \frac{(\Delta x_1)^2}{4\Delta t} \frac{V_{j+1,k}^n - 2V_{j,k}^n + V_{j-1,k}^n}{(\Delta x_1)^2} \\ & \left. - \frac{(\Delta x_2)^2}{4\Delta t} \frac{V_{j,k+1}^n - 2V_{j,k}^n + V_{j,k-1}^n}{(\Delta x_2)^2} \right] \end{aligned}$$

which for  $\Delta x_1 = \Delta x_2 = h$  and  $\lambda = h/\Delta t$  corresponds to a second order approximation of the equation

$$\partial_t v^h + H(\partial_{x_1} v^h, \partial_{x_2} v^h) = \frac{\Delta x^2}{4\Delta t} \sum_i \partial_{x_i x_i} v - \sum_{i,j} \frac{\Delta t}{2} \partial_{x_i} H \partial_{x_j} H \partial_{x_i x_j} v.$$

## Chapter 10

# Rare Events and Reactions in SDE

Transition between stable equilibrium solutions are used to model for instance reaction paths and reaction rates in chemistry and nucleation phenomena in phase transitions excited by thermal fluctuations. An example of such nucleation in an under cooled liquid is the formation of the initial crystal that starts to grow to a whole solid, taking place every year in the first cold calm winter night in Swedish lakes. Deterministic differential equations cannot model such transitions between equilibrium states, since a deterministic solution never escapes from a stable equilibrium. This section shows how stochastic differential equations are used to model reaction paths and its rates, using large deviation theory from an optimal control perspective.

Let us start with a deterministic model

$$\dot{X}^t = -V'(X^t) \quad t > 0,$$

where the potential  $V : \mathbb{R} \rightarrow \mathbb{R}$  is a scalar double well function, see Figure 10.2, with two stable equilibrium points  $x_+$  and  $x_-$ , and one unstable equilibrium point  $x_0$  in between. We see from the phase portrait Figure ?? that

$$\lim_{t \rightarrow \infty} X^t = \begin{cases} x_- & \text{if } X^0 < x_0 \\ x_+ & \text{if } X^0 > x_0 \\ x_0 & \text{if } X^0 = x_0, \end{cases} \quad (10.1)$$

which means that a path from one stable equilibrium point to another stable equilibrium point is not possible in this deterministic setting.

The stochastic setting

$$dX^t = -V'(X^t)dt + \sqrt{2\epsilon}dW^t \quad (10.2)$$

can model transitions between  $x_-$  and  $x_+$ . In this section we focus on the case when the positive parameter  $\epsilon$  (which measures the temperature in the chemistry model) is small, that is we study a small stochastic perturbation of the deterministic case. By

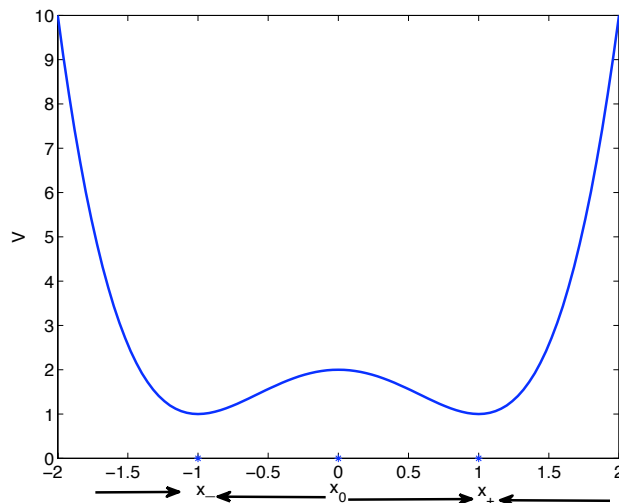


Figure 10.1: Illustration of a double well with two local minima points at  $x_-$  and  $x_+$  and one local maximum point at  $x_0$ .

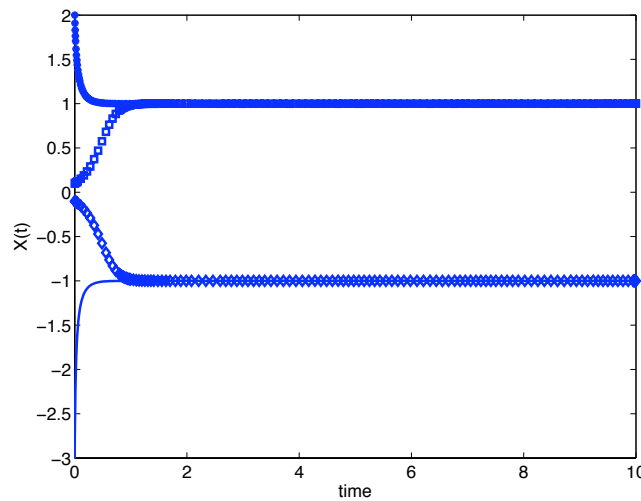


Figure 10.2: Four paths  $X^t$  from a double well potential with two local minima points at  $x_-$  and  $x_+$  and one local maximum point at  $x_0$ .

introducing noise in the model, we may ask what is the probability to jump from one well to the other; since  $\epsilon$  is small these transitions will be rare events. More precisely we shall for the model (10.2) determine:

- the invariant probability distribution and convergence towards it as time tends to infinity,
- the asymptotic behaviour of jumps from one well to another, i.e. reaction rates and reaction paths.

## 10.1 Invariant Measures and Ergodicity

Consider now a stochastic differential equation

$$dX^t = -V'(X^t)dt + \sqrt{2\epsilon}dW^t \quad (10.3)$$

with a potential  $V : \mathbb{R}^d \rightarrow \mathbb{R}$  that is smooth and satisfies  $\int_{\mathbb{R}^d} e^{-V(x)/\epsilon} dx < \infty$ , which implies that  $V(x) \rightarrow \infty$  as  $|x| \rightarrow \infty$ . We also assume a global Lipschitz bound on  $V'$  to have a well defined solution  $X$ , but the global Lipschitz bound can be relaxed. The probability density for an SDE solves the Fokker-Planck equation 4.9. Sometimes this has a time independent solution - the corresponding probability measure is called an *invariant measure*. It is called invariant because if we start with this probability measure as initial probability distribution, the probability distribution obtained from the Fokker-Planck equation for later time remains unchanged, i.e. this probability distribution is time invariant. In the case of an SDE with additive noise and a drift that is the gradient of a potential function, as in (10.3), the invariant measure can be explicitly computed:

**Theorem 10.1.** *The SDE-model (10.3) has the invariant measure*

$$\left( \int_{\mathbb{R}^d} e^{-V(x)/\epsilon} dx \right)^{-1} e^{-V(x)/\epsilon} dx.$$

*Proof.* The Fokker-Planck equation corresponding to the dynamics (10.3) takes the form

$$\partial_t p - \partial_x(V'(x)p(x)) - \epsilon \partial_{xx} p = 0. \quad (10.4)$$

The condition to have an invariant solution means that it is time independent, i.e.  $\partial_t p = 0$ , and the Fokker-Planck equation can be solved explicitly

$$\epsilon p' + V'p = c,$$

for a constant  $c$ . The density  $p$  should be integrable, and consequently  $p(x)$  and  $p'(x)$  must tend to zero as  $|x|$  tends to infinity. Therefore we have  $c = 0$ , which implies

$$\int \frac{dp}{p} = - \int \frac{V'}{\epsilon} dx,$$

with the solution

$$\log p(x) = C' - \frac{V(x)}{\epsilon} \quad \text{for a constant } C',$$

so that for another constant  $C$

$$p(x) = Ce^{-V(x)/\epsilon}.$$

The requirement that  $\int_{\mathbb{R}^d} p(x)dx = 1$  determines the constant to be  $C = \left(\int_{\mathbb{R}^d} e^{-V(x)/\epsilon} dx\right)^{-1}$ .  $\square$

$\square$

A Monte-Carlo method to compute expected values  $\int_{\mathbb{R}^d} g(y)p_0(y)dy$  in an equilibrium environment (with invariant density  $p_0$ ) is typically based on approximations of the integral  $T^{-1} \int_0^T g(X^t)dt$  for large  $T$ ; therefore it is important to understand some basic conditions and properties of such approximations, which is the purpose of the next two theorems.

**Theorem 10.2.** *If one starts with any initial probability density and the density converges time asymptotically to the invariant density  $p_0$ , i.e. for any  $\tau > 0$  the pointwise limit*

$$\lim_{t \rightarrow \infty} \tau^{-1} \int_t^{t+\tau} p^s ds = p^0$$

*is satisfied, then for any continuous bounded function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  there holds in the weak sense*

$$\lim_{T \rightarrow \infty} T^{-1} \int_0^T g(X^t)dt = \int_{\mathbb{R}^d} g(y)p_0(y)dy. \quad (10.5)$$

We say that the stochastic process  $X$  is *ergodic* and that the invariant measure,  $p_0$ , is ergodic if (10.5) holds for all bounded continuous  $g$ .

*Proof.* The proof has two steps - to verify that the expected value converges and then estimate the deviation from this limit.

*Step 1.* By the assumption of the converging density we have

$$\begin{aligned} \lim_{T \rightarrow \infty} \mathbb{E}[T^{-1} \int_0^T g(X^t)dt] &= \lim_{T \rightarrow \infty} \mathbb{E}\left[T^{-1} \left( \int_0^{T^{1/2}} g(X^t)dt + \int_{T^{1/2}}^T g(X^t)dt \right)\right] \\ &= \lim_{T \rightarrow \infty} \mathbb{E}\left[T^{-1} \int_0^{T^{1/2}} g(X^t)dt + T^{-1} \sum_{n=T^{1/2}}^{T-1} \int_n^{n+1} g(X^t)dt\right] \\ &= \underbrace{\int_{\mathbb{R}^d} g(y)p_0(y)dy}_{=:\mathbb{E}_0[g]}, \end{aligned}$$

where the first integral tends to zero, since  $g$  is bounded and  $T^{1/2}/T \rightarrow 0$ , and the  $T - T^{1/2}$  integrals in the sum converge by the assumption, as explained in Example 10.4.



*Step 2.* Let  $T = M\tau$  for some large  $\tau, M$  and write the integral as a sum over  $M$  terms

$$T^{-1} \int_0^T g(X^t) dt = M^{-1} \sum_{n=1}^M \tau^{-1} \int_{n\tau}^{(n+1)\tau} g(X^t) dt.$$

If these terms were independent, the law of large numbers would show that the sum converges almost surely, as  $M$  tends to infinity. Since the terms are only asymptotically independent as  $\tau \rightarrow \infty$ , we need some other method: we shall use Chebyshev's inequality to prove convergence in probability. Let  $\xi_n := \tau^{-1} \int_{n\tau}^{(n+1)\tau} (g(X^t) - \mathbb{E}_0[g]) dt$ , we want to verify that for any  $\gamma > 0$

$$\lim_{M, \tau \rightarrow \infty} P\left(\frac{|\sum_{n=1}^M \xi_n|}{M} > \gamma\right) = 0. \quad (10.6)$$

Chebyshev's inequality implies

$$\begin{aligned} & P\left(\left|\sum_{n=1}^M \xi_n/M\right| > \gamma\right) \\ & \leq \gamma^{-2} \mathbb{E}\left[\sum_n \sum_m \xi_n \xi_m / M^2\right] \\ & = \gamma^{-2} M^{-2} \sum_n \sum_m \tau^{-2} \mathbb{E}\left[\int_n (g(X^t) - \mathbb{E}_0[g]) dt \int_m (g(X^s) - \mathbb{E}_0[g]) ds\right] \\ & = 2\gamma^{-2} M^{-2} \sum_{n>m} \tau^{-2} \int_n \int_m \mathbb{E}\left[\mathbb{E}[(g(X^t) - \mathbb{E}_0[g])(g(X^s) - \mathbb{E}_0[g]) \mid X^s]\right] dt ds \\ & \quad + \gamma^{-2} M^{-2} \sum_n \left(\tau^{-1} \int_n \mathbb{E}[g(X^t) - \mathbb{E}_0[g]] dt\right)^2 \\ & =: I \end{aligned}$$

and since the density  $p^t$  converges we can for each  $\delta > 0$  choose  $\tau$  sufficiently large so that

$$\begin{aligned} I & = 2\gamma^{-2} M^{-2} \sum_{n>m} \tau^{-2} \int_n \int_m \mathbb{E}\left[\int_{\mathbb{R}^d} g(y)(p^t(y) - p_0(y)) dy (g(X^s) - \mathbb{E}_0[g])\right] dt ds \\ & \quad + \gamma^{-2} M^{-2} \sum_n \left(\tau^{-1} \int_n \mathbb{E}[g(X^t) - \mathbb{E}_0[g]] dt\right)^2 \\ & \leq \gamma^{-2} \delta + C\gamma^{-2} M^{-1} \end{aligned}$$

which proves (10.6). □

□

**Theorem 10.3.** *The process  $X$  generated by (10.3) is ergodic for positive  $\epsilon$ .*

*Proof.* Theorem 10.2 tells us that it remains to verify that the probability density converges time asymptotically to the invariant density. Let  $p_0$  be the invariant solution and define the *entropy*

$$E^t := \int_{\mathbb{R}^d} p \log \frac{p}{p_0} dx.$$

We know from Corollary 4.9 that  $p$  is non negative. The proof has three steps: to show that the entropy decays, that the entropy is non negative, and that the decaying entropy implies convergence of the density to the invariant density.

*Step 1.* Show that  $\dot{E}^t = -\epsilon^{-1} \int |\epsilon p' + V'p|^2 p^{-1} dx$ . Differentiation, the Fokker-Planck equation (10.4), and integration by parts<sup>1</sup> imply

$$\begin{aligned} \dot{E}^t &= \int_{\mathbb{R}^d} \partial_t p \log \frac{p}{p_0} + \partial_t p \frac{p}{p_0} dx \\ &= \int_{\mathbb{R}^d} \underbrace{\partial_t p}_{=(V'p)'+\epsilon p''} \left( \log \frac{p}{p_0} + 1 \right) dx \\ &= \int_{\mathbb{R}^d} ((V'p)' + \epsilon p'') \left( \log \frac{p}{p_0} + 1 \right) dx \\ &= - \int_{\mathbb{R}^d} (V'p + \epsilon p') \cdot \left( \frac{p'}{p} - \underbrace{\frac{p_0'}{p_0}}_{-V'/\epsilon} \right) dx \\ &= -\epsilon^{-1} \int_{\mathbb{R}^d} |V'p + \epsilon p'|^2 p^{-1} dx. \end{aligned}$$

*Step 2.* Show that  $E^t \geq 0$  using that  $p$  and  $p_0$  have the same mass and that  $\log x$  is concave. We have

$$E^t = \int_{\mathbb{R}^d} p \log \frac{p}{p_0} dx = \int_{\mathbb{R}^d} p \left( -\log \frac{p_0}{p} + \frac{p_0}{p} - 1 \right) dx$$

and the concavity of the logarithm implies  $\log x \leq x - 1$ , which establishes  $E^t \geq 0$ .

*Step 3.* Time integration of Step 1 gives

$$E^T + \epsilon^{-1} \int_0^T \int_{\mathbb{R}^d} |\epsilon p' + V'p|^2 p^{-1} dx dt = E^0, \quad (10.7)$$

and since  $E^T$  is non negative and  $E^0$  is assumed to be bounded, we see that the integral  $\int_0^T \int |\epsilon p' + V'p|^2 p^{-1} dx dt$  also is bounded uniformly in  $T$ . Therefore we have, for any  $\tau > 0$ , that  $\tau^{-1} \int_t^{t+\tau} \epsilon p'^s + V'p^s ds \rightarrow 0$  in  $L^2(\mathbb{R}^d)$  as  $t \rightarrow \infty$ , which gives  $\tau^{-1} \int_t^{t+\tau} p^s ds \rightarrow p_0$  as follows: integration of

$$\epsilon p'^t + V'p^t =: f^t$$

shows that

$$p(x, t) = e^{-V(x)/\epsilon} \left( C + \int_0^x f(y, t) e^{V(y)/\epsilon} dy \right)$$

---

<sup>1</sup>A better way, in the sense of requiring less assumptions, is to directly study the Fokker-Planck equation in its weak form; then the integration by parts is not needed and (10.7) is obtained directly.

so that  $\tau^{-1} \int_t^{t+\tau} p^s ds \rightarrow p_0$  as  $t \rightarrow \infty$ , since  $\tau^{-1} \int_t^{t+\tau} f^s ds \rightarrow 0$  in  $L^2(\mathbb{R}^d)$ .  $\square$

$\square$

**Example 10.4** (No mass escapes to infinity). The aim here is to verify that the pointwise limit  $\lim_{\tau \rightarrow \infty} \int_{\tau}^{\tau+1} p^t dt = p_0$  implies the weak limit

$$\lim_{\tau \rightarrow \infty} \int_{\tau}^{\tau+1} \int_{\mathbb{R}^d} gp^t dx dt = \int_{\mathbb{R}^d} gp_0 dx, \quad (10.8)$$

for any bounded continuous function  $g$ .

Let  $\bar{p}^{\tau} := \int_{\tau}^{\tau+1} p^t dt$  and define  $\phi : (0, \infty) \rightarrow \mathbb{R}$  by  $\phi(x) = x \log x / p_0$ . The function  $\phi$  is convex and Jensen's inequality implies together with (10.7)

$$\begin{aligned} E^0 &\geq \int_{\mathbb{R}^d} \int_{\tau}^{\tau+1} p^t \log \frac{p^t}{p_0} dt dx \\ &= \int_{\mathbb{R}^d} \int_{\tau}^{\tau+1} \phi(p^t) dt dx \\ &\geq \int_{\mathbb{R}^d} \phi\left(\int_{\tau}^{\tau+1} p^t dt\right) dx \\ &= \int_{\mathbb{R}^d} \phi(\bar{p}^{\tau}) dx. \end{aligned}$$

Therefore we have for any positive number  $n$

$$\int_{\mathbb{R}^d} \bar{p}^{\tau} 1_{\{\bar{p}^{\tau} > np_0\}} dx \leq \frac{E^0}{\log n}. \quad (10.9)$$

We can split our integral into two

$$\int_{\mathbb{R}^d} \int_{\tau}^{\tau+1} gp^t dt dx = \int_{\mathbb{R}^d} g\bar{p}^{\tau} dx = \int_{\mathbb{R}^d} g\bar{p}^{\tau} 1_{\{\bar{p}^{\tau} > np_0\}} dx + \int_{\mathbb{R}^d} g\bar{p}^{\tau} 1_{\{\bar{p}^{\tau} \leq np_0\}} dx,$$

where dominated convergence yields

$$\lim_{\tau \rightarrow \infty} \int_{\mathbb{R}^d} g\bar{p}^{\tau} 1_{\{\bar{p}^{\tau} \leq np_0\}} dx = \int_{\mathbb{R}^d} gp_0 dx$$

and (10.9) shows that the other integral is negligible small

$$\int_{\mathbb{R}^d} g\bar{p}^{\tau} 1_{\{\bar{p}^{\tau} > np_0\}} dx \leq C / \log n$$

as  $n \rightarrow \infty$ , which proves the limit (10.8).

**Exercise 10.5** (Invariant measure for Ornstein-Uhlenbeck). Show that the invariant measure for the Ornstein-Uhlenbeck process is a normal distribution.

**Exercise 10.6** (Vanishing noise density is not the deterministic density). Prove that for a smooth function  $V$  on a bounded set  $A$

$$\lim_{\epsilon \rightarrow 0+} \epsilon \log \int_A e^{-V(y)/\epsilon} dy = - \inf_{y \in A} V(y).$$

Such a limit was first studied by Laplace.

**Exercise 10.7.** Show that for a smooth function  $V$  on a bounded set  $A$  with a unique global minimum point  $y_+$ , the probability density  $\frac{e^{-V(y)/\epsilon}}{\int_A e^{-V(y)/\epsilon} dy}$  has the limit expected value

$$\lim_{\epsilon \rightarrow 0+} \frac{\int_A e^{-V(y)/\epsilon} \phi(y) dy}{\int_A e^{-V(y)/\epsilon} dy} = \phi(y_+),$$

Compare this limit with the time-asymptotic "probability" density for the deterministic  $\epsilon = 0$  case (10.1) and show they are different. What can be concluded about the limits  $t \rightarrow \infty$  and  $\epsilon \rightarrow 0+$  of the probability density?

**Example 10.8** (Simulated Annealing). The stochastic differential equation (10.3) can also be used to find minima of functions  $V : \mathbb{R}^d \rightarrow \mathbb{R}^d$ : we know that its invariant measure has the density  $\frac{\int_A e^{-V(y)/\epsilon} \phi(y) dy}{\int_A e^{-V(y)/\epsilon} dy}$ , which by Exercise 10.7 concentrates at  $x \in \operatorname{argmin} V$ . Therefore, by simulating the stochastic differential equation for very long time with decreasing  $\epsilon$  one expect to have the path  $X$  most of the time in the global minimum; more precisely choose  $\epsilon = \epsilon_1$  for  $t \in [0, T_1], \dots, \epsilon = \epsilon_n$  for  $t \in [T_{n-1}, T_n]$ , with  $\epsilon_n \searrow 0+$  and  $T_n \nearrow \infty$  as  $n \rightarrow \infty$ . This method is called simulated annealing and it can be proven to work for a precise choice of  $\epsilon_n$  and  $T_n$ , see [?]. The advantage with the method is that a global minimum is found and the main question is to find a good combination of  $\epsilon_n$  and  $T_n$  suitable for the particular  $V$  studied.

## 10.2 Reaction Rates

The invariant ergodic measure for  $X$  shows that there is a finite probability to reach all states from any point when  $\epsilon > 0$ , in contrast to the deterministic case  $\epsilon = 0$ ; the invariant measure also shows that these probabilities are exponentially small, proportional to  $e^{-V/\epsilon}$ . It is practical to relate reaction rates to exit times from domains: define for  $X$  solving (10.3) and a given domain  $A \in \mathbb{R}^d$  the exit time

$$\tau(X) = \inf\{t : X^t \notin A\}.$$

We want to understand the exit probability

$$P(\tau < T) = \mathbb{E}[1_{\tau < T}] =: q_\tau \quad \text{as } \epsilon \rightarrow 0+.$$

The Kolmogorov-backward equation shows that

$$\begin{aligned} \partial_t q_\tau - V^l \cdot \partial_x q_\tau + \epsilon \partial_{xx} q_\tau &= 0 & \text{in } A \times (0, T) \\ q_\tau(x, \cdot) &= 1 & \text{on } \partial A \times (0, T) \\ q_\tau(\cdot, T) &= 0 & \text{on } A \times \{T\}. \end{aligned} \tag{10.10}$$

**Remark 10.9** (A useless solution). A naive try could be to remove the diffusion part  $\epsilon \partial_{xx} q_\tau$  in (10.4); that leads to the hyperbolic equation

$$\begin{aligned} \partial_t q_\tau - V' \cdot \partial_x q_\tau &= 0 & \text{in } A \times (0, T) \\ q_\tau &= 1 & \text{on } \partial A \times (0, T) \\ q_\tau(\cdot, T) &= 0 & \text{on } A \times \{T\} \end{aligned} \tag{10.11}$$

which can be solved by the characteristics  $\dot{y}^t = -V'(y^t)$ :

$$\frac{d}{dt} q_\tau(y^t, t) = \partial_t q_\tau + \frac{dy^t}{dt} \cdot \partial_x q_\tau = \partial_t q_\tau - V' \cdot \partial_x q_\tau = 0.$$

Since the equilibrium points are stable, it turns out that all characteristics leave the domain on the upper part  $t = T$  see Figure 10.3, where  $q_\tau = 0$ , so that the solution of (10.11) becomes  $q_\tau = 0$ , and that is a useless solution.

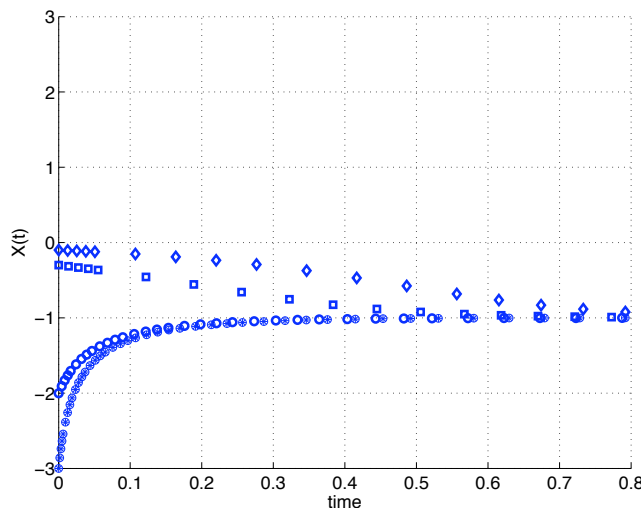


Figure 10.3: Four paths  $X^t$  starting with  $X^0 < x_0$  in the domain of the global attractor  $x_-$

The limit in Remark 10.9 needs to be refined to give something useful. The invariant measure with probabilities proportional to  $e^{-V/\epsilon}$  suggests a change of variables  $q_\tau(x, t) = e^{w_\epsilon(x, t)/\epsilon}$ . The right way to study  $q_\tau$  as  $\epsilon \rightarrow 0+$  is to use the limit

$$\lim_{\epsilon \rightarrow 0+} \epsilon \log q_\tau = \lim_{\epsilon \rightarrow 0+} w_\epsilon =: w$$

which we believe has a bounded non positive limit, using the invariant measure. Since  $q_\tau$  is a probability we know that  $w_\epsilon \leq 0$  and (10.10) implies that  $w_\epsilon$  solves the second order Hamilton-Jacobi equation

$$\begin{aligned} \partial_t w_\epsilon - V' \cdot \partial_x w_\epsilon + \partial_x w_\epsilon \cdot \partial_x w_\epsilon + \epsilon \partial_{xx} w_\epsilon &= 0 & \text{in } A \times (0, T) \\ w_\epsilon(x, \cdot) &= 0 & \text{on } \partial A \times (0, T) \\ w_\epsilon(\cdot, T) &= -\infty & \text{on } A \times \{T\}. \end{aligned}$$

A good way to understand this Hamilton-Jacobi equation is to view it as an optimal control problem. In the limit as  $\epsilon$  tends to zero, the optimal control problem becomes deterministic, see Theorem 9.10; assume that  $\lim_{\epsilon \rightarrow 0^+} w_\epsilon =: w$  to obtain the first order Hamilton-Jacobi equation

$$\begin{aligned} \partial_t w - \underbrace{V' \cdot \partial_x w + \partial_x w \cdot \partial_x w}_{=: H(w(x), x)} &= 0 \quad \text{in } A \times (0, T) \\ w(x, \cdot) &= 0 \quad \text{on } \partial A \times (0, T) \\ w(\cdot, T) &= -\infty \quad \text{on } A \times \{T\}. \end{aligned}$$

Following Section 9.1.4, a useful optimal control formulation for this Hamilton-Jacobi equation is

$$\begin{aligned} \dot{Y}^t &= -V'(Y^t) + 2\alpha^t \\ \max_{\alpha: (0, T) \rightarrow \mathbb{R}^d} & - \int_0^\tau |\alpha^t|^2 dt + g(Y^\tau, \tau) \end{aligned}$$

which has the right Hamiltonian

$$\sup_{\alpha \in \mathbb{R}^d} \left( \lambda \cdot (-V'(y) + 2\alpha) - |\alpha|^2 \right) = H(\lambda, y) = -V'(y) \cdot \lambda + |\lambda|^2.$$

Here the final cost is zero, if the exit is on the boundary  $\partial A \times (0, T)$ , and minus infinity if the exit is on  $A \times \{T\}$  (i.e. the path did not exit from  $A$ ):

$$g(x, t) = \begin{cases} 0 & \text{on } \partial A \times (0, T) \\ -\infty & \text{on } A \times \{T\}. \end{cases}$$

Theorem 9.10 shows that the limit  $\lim_{\epsilon \rightarrow 0^+} \epsilon \log q_\tau = \lim_{\epsilon \rightarrow 0^+} w_\epsilon = w$  satisfies

$$\begin{aligned} w(x, t) &= \sup_{\alpha: (t, \tau) \rightarrow \mathbb{R}^d} - \int_t^\tau |\alpha|^2 dt + g(Y^\tau, \tau) \\ &= \sup_{\alpha} - \frac{1}{4} \int_t^\tau |\dot{Y}^t + V'(Y^t)|^2 dt + g(Y^\tau, \tau). \end{aligned}$$

When  $T$  tends to infinity and  $X^0$  is an equilibrium point, this limit  $w$  has a simple explicit solution showing that reaction rates are determined from local minima and saddle points of  $V$ , cf. Figure 10.4:

**Theorem 10.10.** *Assume that  $y_+$  is a global attractive equilibrium in  $A$ . Let  $X^0 = y_+$ , then*

$$\lim_{T \rightarrow \infty} \lim_{\epsilon \rightarrow 0^+} \epsilon \log q_\tau = V(y_+) - \inf_{y \in \partial A} V(y). \quad (10.12)$$

*Proof.* It is clear the optimal control paths starting in  $y_+$  need to exit through  $\partial A$ , so  $g(Y^\tau) = 0$ . The integral cost can be rewritten as

$$\begin{aligned} & \sup_{\alpha} -\frac{1}{4} \int_0^\tau |\dot{Y}^t + V'(Y^t)|^2 dt \\ &= \sup_{\alpha} \left( -\frac{1}{4} \int_0^\tau \underbrace{|\dot{Y}^t - V'(Y^t)|^2}_{\geq 0} dt - \underbrace{\int_0^\tau \dot{Y}^t \cdot V'(Y^t) dt}_{V(Y^\tau) - V(y_+)} \right). \end{aligned} \quad (10.13)$$

Here the last integral is minimal if  $Y^\tau$  exits through a point on  $\partial A$  where  $V$  is minimal, which is a *saddle point* if we have chose  $A$  to be the largest domain where  $y_+$  is a global attractor. It remains to show that such an exit is compatible with having the first integral equal to zero; the first integral equals zero means that  $\dot{Y}^t = V'(Y^t)$ , which implies that  $Y$  moves orthogonal to the level lines of the  $V$ -potential. Such a path is possible by taking  $\alpha = V'(Y^t)$  and requires  $T$  to be sufficiently large so that the time to reach the boundary on the optimal path  $\dot{Y}^t = V'(Y^t)$  is shorter, when  $X^0$  tends to  $y_+$  this time tends to infinity.  $\square$

$\square$

We see that the probability to exit from an equilibrium is exponentially tiny, proportional to  $e^{-(\inf_{y \in \partial A} V(y) - V(y_+))/\epsilon}$  as  $\epsilon$  tends to zero, and therefore such exits are rare events. In the next section we show that the most probable path, the so called *reaction paths*, that gives such rare events are those where the stochastic paths  $X$  closely follow the optimal control paths  $Y$ . Since  $\epsilon$  is small and the control  $\alpha$  is not, the Brownian motion must some time be large of order  $\epsilon^{-1/2}$ . Therefore the rare events of exits depend on the rare events of such large deviation in the Brownian motion.

The Theorem relates to the basis of reaction theory in chemistry and statistical physics, where the probability to go from one state with energy  $V_1$  to another with energy  $V_2 > V_1$  is proportional to Boltzmanns rate  $e^{-(V_2 - V_1)/(k_B \mathcal{T})}$ ; here  $k_B$  is Boltzmanns constant and  $\mathcal{T}$  is the temperature. We see that, with  $\epsilon = k_B \mathcal{T}$  and  $V$  the energy, the simple model (10.3) can describe reactions and physical transition phenomena. A simple way to see that the *reaction rate* is  $q_\tau$  is to take  $N$  independent particles starting in  $y_+$ . After very long time  $Nq_\tau$  of them have exited from the domain and the reaction rate becomes the quotient  $Nq_\tau/N = q_\tau$ .

**Exercise 10.11.** Show that the mean exit time  $u_\epsilon(x, t) := \mathbb{E}[\tau - t \mid X^t = x]$  satisfies

$$\lim_{\epsilon \rightarrow 0^+} \epsilon \log u_\epsilon(y_+, t) = \inf_{y \in \partial A} V(y) - V(y_+).$$

**Exercise 10.12.** Does

$$\lim_{\epsilon \rightarrow 0^+} \epsilon \log q_\tau = V(X^0) - \inf_{y \in \partial A} V(y)$$

hold when  $X^0$  starts from a different point than the global attractor in  $A$ ? Answer: sometimes but not in general depending on  $X^0$  - how?

Exercise 10.11 shows that the product of the limits of the mean exit time and the probability to exit is equal to one, that is the mean exit time is exponentially large, roughly  $e^{(\inf_{y \in \partial A} V(y) - V(y_+))/\epsilon}$ .

### 10.3 Reaction Paths

This section motivates why the most probable exit paths  $X$  closely follow the optimal control paths  $Y$ . We saw in Theorem 10.10 that in the case  $T$  tending to infinity and  $Y^0 = y_+$ , the optimal path  $Y$  is orthogonal to the level sets of the potential  $V$  and the path starts from the minimum point  $y_+$  (where  $V(y_+) = \min_{y \in A} V(y)$ ) and moves towards the minimum on the boundary  $\operatorname{argmin}_{y \in \partial A} V(y)$ , see Figure 10.4. For bounded  $T$  the situation may change and the time to reach the boundary with the control  $\alpha = V'$  may be larger than  $T$ , so that the first integral in (10.13) does not vanish and the optimal control becomes different; therefore also the exit probability is different and (10.12) is invalid; clearly such early time exit probabilities are also interesting when a rare event is unwanted, e.g. for hard-disc and power-plant failures. These most probable paths following the optimal control paths are called the *reaction paths*. Since the exit probability is small and the most probable exit path makes a *large deviation* from the equilibrium on a time span of order one, which is small compared to the expected exit time of order  $e^{C/\epsilon}$  (for some positive  $C$ ), the exit process can on long time spans be considered as a Poisson process with the rate  $1/\mathbb{E}[\tau - t] \simeq q_\tau$ .

To verify that the most probable exit paths follow the optimal control paths, we want to in some sense relate the stochastic increments  $\sqrt{2\epsilon} dW^t$  with the control increments  $\alpha^t dt$ . Our first step in this direction is to find a probability measure on whole paths  $X$ , and then to see how probable the  $X$ -paths close to the optimal control paths  $Y_*$  are compared to the  $X$ -paths away from  $Y_*$ . It is clear that the probability to find  $X = Y_*$  is zero, so we need to modify this argument somewhat. An informal way to understand the probability of whole paths is to consider Euler discretizations of (10.3)

$$\left(\frac{\Delta X}{\Delta t} + V'(X_i)\right)\Delta t = \sqrt{2\epsilon}\Delta W$$

with the probability density

$$\begin{aligned} P(\Delta W = y_i) &= e^{-\frac{|y_i|^2}{2\Delta t}} \frac{dy_i}{(2\pi\Delta t)^{d/2}} \\ &= e^{-|\frac{\Delta X}{\Delta t} + V'(X_i)|^2 \Delta t / (4\epsilon)} \frac{dy_i}{(2\pi\Delta t)^{d/2}}. \end{aligned}$$

Therefore the probability measure for a whole path is

$$\begin{aligned} \prod_{i=1}^n e^{-\frac{|y_i|^2}{2\Delta t}} \frac{dy_i}{(2\pi\Delta t)^{d/2}} &= \prod_{i=1}^n e^{-|\frac{\Delta X}{\Delta t} + V'(X_i)|^2 \Delta t / (4\epsilon)} \frac{dy_i}{(2\pi\Delta t)^{d/2}} \\ &= e^{-\sum_{i=1}^n |\frac{\Delta X}{\Delta t} + V'(X_i)|^2 \Delta t / (4\epsilon)} \frac{dy_1}{(2\pi\Delta t)^{d/2}} \cdots \frac{dy_n}{(2\pi\Delta t)^{d/2}}. \end{aligned}$$



The most probable path is the one that maximises the probability density

$$e^{-\sum_{i=1}^n |\frac{\Delta X}{\Delta t} + V'(X_i)|^2 \Delta t / (4\epsilon)},$$

this is called the *maximum likelihood method*. In the previous section we saw that the optimal control problem does precisely this maximisation. Therefore the optimal control paths generate the most probable stochastic paths. If the density in the maximum likelihood method is almost uniform, the result is doubtful. Here the situation is the opposite - when  $\epsilon$  tends to zero, the density concentrates on the most probable event, see Exercise 10.13.

If we consider  $W$  or  $\alpha$  as perturbations, we see that the solution we have obtained is the solution of the *least-squares problems*

$$\min_W \int_0^\tau |\dot{X}^t + V'(X^t)|^2 dt = \min_\alpha \int_0^\tau |\dot{Y}^t + V'(Y^t)|^2 dt,$$

where  $\dot{X}^t + V'(X^t)$  and  $\dot{Y}^t + V'(Y^t)$  are the residuals, that is the error in the equation.

**Exercise 10.13.** In the limit as  $\epsilon$  tends to zero, we saw in Exercise 10.6 that if  $\int_A e^{-V(y)} dy$  is bounded, then

$$\lim_{\epsilon \rightarrow 0^+} \epsilon \log \int_A e^{-V(y)/\epsilon} dy = - \inf_{y \in A} V(y).$$

Show that for a smooth function  $f$  on a bounded set  $A$  with a unique maximum point  $y_+$ , the probability density  $\frac{e^{f(y)/\epsilon}}{\int_A e^{f(y)/\epsilon} dy}$  has the limit expected value

$$\lim_{\epsilon \rightarrow 0^+} \frac{\int_A e^{f(y)/\epsilon} \phi(y) dy}{\int_A e^{f(y)/\epsilon} dy} = \phi(y_+),$$

which means that in the limit the most probable event almost surely happens and nothing else.

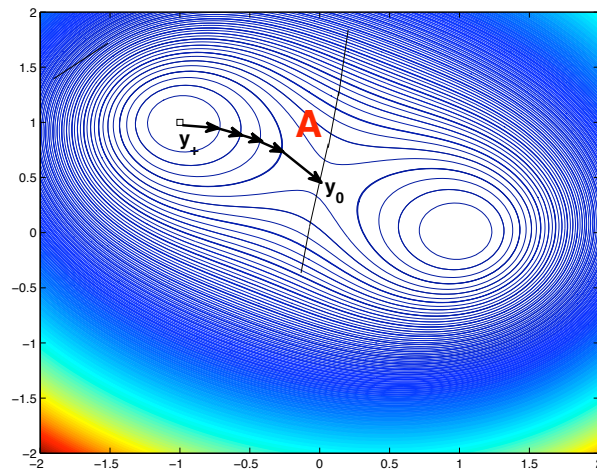


Figure 10.4: The optimal reaction path starting in the attractor  $y_+$  moving to the saddlepoint  $y_0 = \operatorname{argmin}_{y \in \partial A}(V(y))$ , inside the domain  $A$  to the left of the dashed line.

## Chapter 11

# Molecular dynamics

The starting point for modelling molecular systems is the eigenvalue problem of the *time-independent Schrödinger equation*

$$H\Psi = E\Psi$$

where the unknown eigenvector  $\Psi$  is a complex valued wave function, depending on the variables of coordinates and spins of all,  $M$ , nuclei and,  $N$ , electrons in the problem, the real number  $E$  is the unknown eigenvalue, and  $H$  is the given Hamiltonian Hermitian operator precisely defined by well known fundamental constants of nature and the Coulomb interaction of all nuclei and electrons. An important issue is its high computational complexity for problems with more than a few nuclei, due to the high dimension of  $\Psi$  which is roughly in  $L^2(\mathbb{R}^{3(M+N)})$ , see [CL]. Already simulation of a single water molecule requires a partial differential equation in 39 space dimensions, which is a demanding task to solve. Therefore coarse-grained approximations are often necessary. The next sections describe the following five useful levels of coarse-grained models:

- In quantum classical molecular dynamics, also called *Ehrenfest dynamics*, the nuclei dynamics is approximated by classical paths, which introduces time and the *time-dependent Schrödinger equation* for the electron dynamics.
- In the *Born-Oppenheimer approximation* the electron wave function in the Ehrenfest dynamics is approximated by the electron ground state for the current nuclei position. This Born-Oppenheimer approximation leads to a molecular system described by a Hamiltonian system, which simulates an ensemble with constant number of particles, volume and energy  $M\ddot{X}_t = -V'(X_t)$ .
- In a situation where one instead wants to simulate a system with constant number of particles, volume and temperature  $T$ , the Born-Oppenheimer approximation can be refined, by including a perturbation of the ground state; this leads to stochastic *Langevin dynamics*  $Mdv_t = -(V'(X_t) + v_t/\tau)dt + \sqrt{T/\tau}dW_t$ . The Langevin dynamics introduces a friction parameter  $1/\tau$ .

- In the high friction limit,  $\tau \rightarrow 0+$ , of Langevin dynamics for long time, the velocity variable  $\dot{X} = v_t$  can be eliminated and the nuclei positions  $X_{s/\tau} \rightarrow \bar{X}_s$  satisfy the *Smoluchowski dynamics*  $d\bar{X}_s = -V'(\bar{X}_s) + \sqrt{T}dW_s$ .
- The next step in the coarse-graining process is to derive partial differential equations for the mass, momentum and energy of a continuum fluid from Langevin or Smoluchowski molecular dynamics, which determines the otherwise unspecified pressure, viscosity and heat conductivity; we present a derivation related to the work by Irvine & Kirkwood (1950) and Hardy (1981).

## 11.1 Molecular dynamics at constant temperature: Zwanzig's model and derivation of Langevin dynamics

This section reviews the Hamiltonian system heat bath model of Zwanzig [Zwa73], with his derivation of stochastic Langevin dynamics, related to the earlier work [FKM65]. Here the model is heavy particles interacting with many light particles – modelling heavy particles in a heat bath of light particles. The model is as simple as possible to have the desired qualitative properties of a system interacting with a heat bath, the following sections then applies a similar formulation to a more fundamental model for nuclei electron systems. The goal here is to give some understanding of simulating, at constant temperature, the coarse-grained molecular dynamics of the heavy particle without resolving the lighter particles, using Langevin dynamics. It is an example how stochastics enter into a coarse-grained model through elimination of some degrees of freedom in a deterministic model, described by a Hamiltonian system. The original model is time reversible while the coarse-grained model is not.

We study  $N_h$  heavy particles and consider particle's positions  $X \in \mathbb{R}^{3N_h}$  in a heat bath with several light particles positioned in  $y_n \in \mathbb{R}^3$ ,  $n = 1, \dots, N$ , relative to the individual equilibrium positions corresponding to  $y_n = 0$ . The harmonic interaction potential

$$U(X) + \sum_{n=1}^N \frac{m\omega_n^2}{2} \left| y_n - \frac{\gamma_n \cdot X}{\omega_n^2} \right|^2,$$

yields the Hamiltonian

$$H := U(X) + \sum_{n=1}^N \frac{m\omega_n^2}{2} \left| y_n - \frac{\gamma_n \cdot X}{\omega_n^2} \right|^2 + \frac{M|\dot{X}|^2}{2} + \sum_n \frac{m|\dot{y}_n|^2}{2} \quad (11.1)$$

and the dynamics

$$M\ddot{X}_t = -U'(X_t) + \sum_n m\omega_n^2 \left( y_n(t) - \frac{\gamma_n \cdot X_t}{\omega_n^2} \right) \frac{\gamma_n}{\omega_n^2}, \quad (11.2)$$

$$m\ddot{y}_n(t) = -m\omega_n^2 \left( y_n(t) - \frac{\gamma_n \cdot X_t}{\omega_n^2} \right). \quad (11.3)$$

Here  $m$  and  $M$  are the light and heavy particle masses, respectively, the function  $U$  is the potential for external forces on the heavy particle and  $\omega_n$  is the particle frequency of oscillation of the light particle,  $n$ , and  $\gamma_n \in \mathbb{R}^{3N_h}$  measures its coupling to heavy particles. Given the path  $X$ , the linear equation (11.3) can be solved explicitly, e.g. with Laplace transform, with the solution

$$y_n(t) - \frac{\gamma_n}{\omega_n^2} X(t) = \underbrace{\sqrt{\frac{k_B \mathcal{T}}{m \omega_n^2}} (\alpha_n \sin(\omega_n t) + \beta_n \cos(\omega_n t))}_{z_n(t)} - \frac{1}{\omega_n^2} \int_0^t \gamma_n \cdot \dot{X}(t-s) \cos(\omega_n s) ds. \quad (11.4)$$

Let both the initial position and velocity for the heavy particle be zero. We assume that the many initial positions and velocities of the light particles are impossible to measure and determine precisely. Clearly, to predict the dynamics of the heavy particle some information of the light particle initial data is necessary: we shall use the equilibrium probability distribution for the light particles depending only on one parameter – the temperature. Section 11.2 presents a motivation of the stochastic model where the initial positions and velocities for the light particles are randomly sampled with the Gibbs probability measure

$$\begin{aligned} & Z^{-1} \exp(-H(y, \dot{y})/(k_B \mathcal{T})) dy_1 \dots dy_N d\dot{y}_1 \dots d\dot{y}_N, \\ Z & := \int_{\mathbb{R}^{2N}} \exp(-H(y, \dot{y})/(k_B \mathcal{T})) dy_1 \dots dy_N d\dot{y}_1 \dots d\dot{y}_N, \\ H(X, \dot{X}, y, \dot{y}) & := U(X) + \sum_{n=1}^N \frac{m \omega_n^2}{2} |y_n - \frac{\gamma_n \cdot X}{\omega_n^2}|^2 + \frac{M |\dot{X}|^2}{2} + \sum_n \frac{m |\dot{y}_n|^2}{2}, \end{aligned} \quad (11.5)$$

which generates  $\alpha_n \in \mathbb{R}^3$  and  $\beta_n \in \mathbb{R}^3$  to be independent stochastic variables with independent standard normal distributed components with zero mean and variance 1.

Inserted into the equation (11.2), for the heavy particle, the solution (11.4) implies that

$$M \ddot{X}_t = -U'(X_t) - \int_0^t \sum_n \frac{m \gamma_n^2}{\omega_n^2} \cos(\omega_n s) \dot{X}(t-s) ds + \underbrace{\sum_n m z_n(t) \gamma_n}_{\zeta(t)} \quad (11.6)$$

where the covariance of the Gaussian process,  $\zeta : [0, \infty) \times \{\text{probability outcomes}\} \rightarrow \mathbb{R}^3$ ,

$$\begin{aligned} \mathbb{E}[\zeta_s^i \zeta_t^i] &= k_B \mathcal{T} \sum_{n=1}^N \frac{m (\gamma_n^i)^2}{\omega_n^2} \cos \omega_n(t-s) =: k_B \mathcal{T} f(t-s), \\ \mathbb{E}[\zeta_s^i \zeta_t^j] &= 0, \quad i \neq j, \end{aligned}$$

also is the integral kernel for the friction term in the generalized Langevin equation (11.6), forming a version of Einstein's fluctuation-dissipation result.

Assume now that the harmonic oscillators are distributed so that the sum over particles is in fact an integral over frequencies with a Debye distribution, i.e. for any function  $h$

$$N^{-1} \sum_{n=1}^N h(\omega_n) \rightarrow \int_0^{\omega_d} h(\omega) \frac{3\omega^2}{\omega_d^3} d\omega,$$

and let  $\gamma_n = \gamma N^{-1/2}$  to obtain

$$M^{-1} f(t) = \frac{3m(\gamma^i)^2}{M\omega_d^3} \frac{\sin \omega_d t}{t},$$

which formally leads to the Langevin equation

$$\begin{aligned} dX_t &= v_t dt, \\ dv_t &= \left( -M^{-1}U'(X) - \tau^{-1}v_t \right) dt + \sqrt{\frac{2k_B\mathcal{T}}{\tau M}} dW_t, \end{aligned} \tag{11.7}$$

as  $\omega_d \rightarrow \infty$  and  $\frac{3\pi m(\gamma^i)^2}{2M\omega_d^3} \rightarrow \tau^{-1}$ , where  $W$  is the standard Wiener process with independent components in  $\mathbb{R}^3$ . This Langevin equation has the unique invariant probability density (that is the time independent solution of the corresponding Kolmogorov forward equation)

$$\frac{e^{-(M|\dot{X}|^2/2+U(X))/\mathcal{T}} dX d\dot{X}}{\int_{\mathbb{R}^6} e^{-(M|\dot{X}|^2/2+U(X))/\mathcal{T}} dX d\dot{X}},$$

which is the heavy particle marginal distribution of the Gibbs distribution

$$\frac{e^{-H(X,\dot{X},y,\dot{y})/\mathcal{T}} dX d\dot{X} dy d\dot{y}}{\int_{\mathbb{R}^{6(N+1)}} e^{-H(X,\dot{X},y,\dot{y})/\mathcal{T}} dX d\dot{X} dy d\dot{y}}$$

in (11.5). We conclude that sampling the light particles from the light particle marginal of the Gibbs distribution leads time asymptotically to having the heavy particle in the heavy particle marginal of the Gibbs distribution: this fundamental stability and consistency property is unique to the Gibbs distribution, as explained in the next section.

Sections ?? to ?? derive *ab initio* Langevin dynamics for nuclei from the Schrödinger equation of interacting nuclei and electrons, in a spirit inspired by Zwanzig's derivation above but using consistency of value functions instead of explicit solutions. The idea of error analysis with value functions is sketched in Section 1.5 and also used for weak convergence of approximations of stochastic differential equations in Section 5.2 and for approximation of optimal control problems in Sections 9.1.7 and 9.2.3.

## 11.2 The Gibbs distribution derived from dynamic stability

At the heart of Statistical Mechanics is the Gibbs distribution

$$\frac{e^{-H(Y,Q)/\mathcal{T}} dY dQ}{\int_{\mathbb{R}^{6N}} e^{-H(Y,Q)/\mathcal{T}} dY dQ}$$

for an equilibrium probability distribution of a Hamiltonian dynamical system

$$\begin{aligned}\dot{Y}_t &= \partial_Q H(Y_t, Q_t) \\ \dot{Q}_t &= -\partial_Y H(Y_t, Q_t)\end{aligned}\tag{11.8}$$

in the canonical ensemble of constant number of particles  $N$ , volume and temperature  $T$ . Every book on Statistical Mechanics gives a motivation of the Gibbs distribution, often based on entropy considerations, cf. [Fey98]. Here we motivate the Gibbs distribution instead from dynamic stability reasons. Consider a Hamiltonian system with light and heavy particles, with position  $Y = (X, y)$ , momentum  $Q = (P, q)$  and the Hamiltonian  $H = H_1(X, P) + H_2(X, y, q)$ , as in (11.1). Assume that it is impractical or impossible to measure and determine the initial data for the light particles. Clearly it is necessary to give some information on the data to determine the solution at a later time. In the case of molecular dynamics it is often sufficient to know the distribution of the particles to determine thermodynamic relevant properties, as e.g. the pressure-law. We saw in Section 11.1 that if the light particles have an initial probability distribution corresponding the Gibbs distribution conditioned on the heavy particle, then the invariant distribution for the heavy particle is unique (in the limit of the Langevin equation) and given by the Gibbs marginal distribution for the heavy particle

$$\frac{e^{-H_1(X,P)/T} dX dP}{\int_{\mathbb{R}^{6N_h}} e^{-H_1(X,P)/T} dX dP}.$$

This stability that an equilibrium distribution of light particles leads to the marginal distribution of the heavy particles holds only for the Gibbs distribution in the sense we shall verify below. This is a desired stability and consistency result:

- (C) we start from an equilibrium density and consider the dynamics of the heavy particles, with the light particles initially distributed according to the light particle equilibrium distribution conditioned on the heavy particles, and end up after long time with the heavy particles distributed according to the heavy particle marginal of the original equilibrium measure; consequently the behavior after long time is consistent with the assumption to start the light particles with this particular equilibrium distribution.

It is in fact this uniqueness of the Gibbs initial probability distribution that makes a stochastic model of the dynamics useful: if we would have to seek the initial distribution among a family of many distributions we could not predict the dynamics in a reasonable way.

To derive this uniqueness of the Gibbs density, we consider first all equilibrium densities of the the Hamiltonian dynamics and then use the consistency check (C) of an equilibrium density and its light particle equilibrium distribution leading to the heavy particle marginal equilibrium distributions to rule out all except the Gibbs density. There are many equilibrium distributions for a Hamiltonian system: the Liouville equation (i.e. the Fokker-Planck equation in the case of zero diffusion)

$$\underbrace{\partial_t f(H)}_{=0} + \partial_Y(\partial_Q H f(H)) - \partial_Q(\partial_Y H f(H)) = 0$$

shows that any positive function  $f$ , depending only on the Hamiltonian  $H$  and not on time, is an invariant probability distribution

$$\frac{f(H(Y, Q)dYdQ}{\int_{\mathbb{R}^{6N}} f(HY, Q)dYdQ}$$

for the Hamiltonian system (11.8). There may be other invariant solutions which are not functions of the Hamiltonian but these are not considered here. Our basic question is now – which of these functions  $f$  have the fundamental property that their light particle distribution generates a unique invariant measure given by the heavy particle marginal distribution? We have seen that the Gibbs distribution is such a solution. Are there other?

Write  $H = H_1 + H_2$  and assume that the number of heavy particles  $N_h$  dominates the number of light particles  $N$ . Then we have

$$\frac{H_2}{H_1} = \mathcal{O}\left(\frac{N}{N_h}\right) \ll 1. \quad (11.9)$$

Let

$$-\log f(H) = g(H)$$

and consider perturbations of the Gibbs distribution in the sense that the function  $g$  satisfies for a constant  $C$

$$\begin{aligned} \lim_{H \rightarrow \infty} \frac{g''(H)H}{g'(H)} &\leq C \\ \lim_{H \rightarrow \infty} \frac{g'(H)H}{g(H)} &\leq C \end{aligned} \quad (11.10)$$

for instance, any monomial  $g$  satisfies (11.10). Taylor expansion yields for some  $\alpha \in (0, 1)$

$$\begin{aligned} -\log f(H) &= g(H_1 + H_2) \\ &= g(H_1) + H_2(g'(H_1) + 2^{-1}g''(H_1 + \alpha H_2)H_2) \end{aligned}$$

and (11.9) and (11.10) implies the leading order term

$$-\log f(H) \simeq g(H_1) + H_2g'(H_1).$$

Define the constant  $T = 1/g'(H_1(X_0, P_0))$ ; the light particle distribution is then asymptotically given by

$$\frac{e^{-H_2/T} dydq}{\int e^{-H_2/T} dydq}.$$

This initial distribution corresponds to a Gibbs distribution with the temperature  $T = 1/g'(H_1(X_0, P_0))$  and the derivation of (11.7) leads to the heavy particle equilibrium distribution

$$\frac{e^{-H_1/T} dXdP}{\int e^{-H_1/T} dXdP}. \quad (11.11)$$



The equilibrium density  $f$  has by (11.9) and (11.10) the leading order expansion

$$\begin{aligned} -\log f(H) &= g(H_1 + H_2) \\ &= g(H_1) + g'(H_1 + \alpha H_2)H_2 \\ &\simeq g(H_1), \end{aligned}$$

which leads to the heavy particle marginal distribution

$$\frac{e^{-g(H_1)} dX dP}{\int e^{-g(H_1)} dX dP} \quad (11.12)$$

The consistency requirement to have the heavy particle distribution (11.11) equal to the heavy particle marginal distribution (11.12) implies that

$$g(H_1) = H_1/T.$$

We conclude that the quotient  $-H/\log f(H)$  is constant, where  $-H/\log f(H) = T$  is called the temperature, and we have derived the Gibbs density  $f(H) = e^{-H/T}$ .

### 11.3 Smoluchowski dynamics derived from Langevin dynamics

See Section 6 in "A stochastic phase-field model derived from molecular dynamics" on <http://www.nada.kth.se/~szepessy/papers.html>.

### 11.4 Macroscopic conservation laws for compressible fluids motivated from molecular dynamics

Molecular dynamics can be used to determine properties of bulk in addition to observables related to smaller nuclei-electron systems. In this section we study the continuum limit of molecular dynamics, which gives us the important connection between microscopic molecular dynamics variables and macroscopic bulk properties as the density, stress, velocity and their conservation of mass, momentum and energy. In particular we will see that a complete macroscopic description of a compressible fluid requires a constitutive relation determining the stress tensor as a function of the density, velocity and energy, which is based on microscopic quantum mechanics. The derivation<sup>1</sup> also gives some insight to simulating molecular dynamics in the different ensembles of  $NVT$  and  $NPT$ , with constant number of particles, volume and temperature, respectively constant number of particles, pressure and temperature.

For a given constant mean velocity  $u$ , the Langevin equation can (with the change of variables  $X_j^t$  replaced by  $X_j^t + tu$ ) for the case with a pair potential be written

$$\begin{aligned} \dot{X}_j^t &= p_j^t - u \\ \dot{p}_j^t &= -\sum_{i \neq j} \Phi'(X_j^t - X_i^t) - K p_j^t + (2KT)^{1/2} \dot{W}_j^t, \end{aligned} \quad (11.13)$$

---

<sup>1</sup> previous related work by Irving & Kirkwood (1950) and Hardy (1981)

and its unique invariant measure is the Gibbs measure

$$\frac{e^{-\left(\sum_j |p_j - u|^2/2 + \sum_j \sum_{i \neq j} \Phi(X_j - X_i)/2\right)/T} dX dp}{\int e^{-\left(\sum_j |p_j - u|^2/2 + \sum_j \sum_{i \neq j} \Phi(X_j - X_i)/2\right)/T} dX dp}$$

for any constant positive temperature  $T$ . We shall study the limit  $K \rightarrow 0+$  as the friction vanishes and then we obtain a Hamiltonian system. To study the continuum limit of molecular dynamics, we consider subsets  $B$  of particles and split the force into

$$\sum_{i \in B, i \neq j} \Phi'(X_j - X_i) + \sum_{i \in B^c} \Phi'(X_j - X_i),$$

where the last sum is the external force due to particles outside  $B$  interacting with particle  $j$  in  $B$  and the notation  $B^c$  means the complement set of  $B$ . To formulate a Hamiltonian for the dynamics of particle in such a set  $B$ , we introduce the characteristic paths  $y^t$  by  $y^t = tu$  and an additional non interacting particle, whose position  $X_0$  measures time  $t$ . We also consider the external potential  $R : \mathbb{R}^3 \times \mathbb{R}_+ \rightarrow \mathbb{R}$  defined by

$$\sum_{i \in B^c} \Phi(X_j^t - X_i^t) = \sum_{i \in B^c} \Phi(X_j^t + y^t - X_i^t - y^t) =: 2R(y^t + X_j^t, t)$$

as a given function of the internal positions  $X_j$  for  $j \in B$ . The local Hamiltonian energy given by

$$\hat{H} := \frac{1}{2} \sum_{j \in B} |p_j^t - u|^2 + \underbrace{\frac{1}{2} \sum_{j \in B} \sum_{i \in B, i \neq j} \Phi(X_j^t - X_i^t)}_{=: V(X)} + \sum_{j \in B} 2R(y^t + X_j^t, t) + p_0^t$$

then yields the vanishing friction dynamics of (11.13) for  $j \in B$

$$\begin{aligned} \dot{X}_j^t &= p_j^t - u \\ \dot{p}_j^t &= - \sum_{i \in B, i \neq j} \Phi'(X_j^t - X_i^t) - \underbrace{\sum_{i \in B^c} \Phi'(X_j^t - X_i^t)}_{=: 2R'(y^t + X_j^t, t)} \end{aligned} \quad (11.14)$$

and  $\dot{X}_0 = 1$ , so that  $X_0 \equiv t$ . Define also the local equilibrium energy function

$$\begin{aligned} H_B &:= \frac{1}{2} \sum_{j \in B} |p_j^t - u|^2 + \frac{1}{2} \sum_{j \in B} \sum_{i \in B, i \neq j} \Phi(X_j^t - X_i^t) + \frac{1}{2} \sum_{j \in B} \sum_{i \in B^c} \Phi(X_j^t - X_i^t) \\ &= \frac{1}{2} \sum_{j \in B} |p_j^t - u|^2 + \frac{1}{2} \sum_{j \in B} \sum_{i \in B, i \neq j} \Phi(X_j^t - X_i^t) + \sum_{j \in B} R(y^t + X_j^t, t) \\ &= \hat{H} - \frac{1}{2} \sum_{j \in B} \sum_{i \in B^c} \Phi(X_j^t - X_i^t) \\ &= \hat{H} - \sum_{j \in B} R(y^t + X_j^t, t) \end{aligned}$$

corresponding to the energy terms in the Gibbs measure related to the particles in  $B$  and note that half the external field is the difference between the local energy  $\hat{H}$  and the local Gibbs energy  $H_B$ . If the external field  $R(\cdot, \cdot)$  is given, the Langevin dynamics (11.13), for  $j \in B$ , has the unique invariant measure

$$\frac{e^{-H_B(X,p)/T} dX dp}{\int e^{-H_B(X,p)/T} dX dp},$$

where the components of  $X$  and  $p$  are restricted to the set  $B$  and  $X_0 = t$ . Theorem 10.3 shows convergence towards the unique invariant measure for the Smoluchowski molecular dynamics, corresponding to the high friction limit, when the given potential is such that the Gibbs measure is bounded. This measure is also one of the invariant measures for the vanishing friction limit (11.14) but to have convergence towards a unique one we consider the Hamiltonian dynamics (11.14) with a vanishing friction parameter  $K$ . The challenge to obtain a complete proof of the continuum limit, would require to verify that convergence to local equilibrium takes place also when the external field  $R$  is not a given function but determined from the almost equilibrium dynamics of the neighboring sets  $B$ .

The convergence towards local equilibrium motivates that also in a case when the mean velocity  $u : \mathbb{R}^3 \times \mathbb{R}_+ \rightarrow \mathbb{R}^3$  and the temperature  $T : \mathbb{R}^3 \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$  are differentiable functions varying on the macroscopic scale, the microscopic set of particles  $B$  see almost constant  $u$  and  $T$ , so that the dynamics relax to the local equilibrium  $\frac{e^{-H_B(X,p)/T} dX dp}{\int e^{-H_B(X,p)/T} dX dp}$  on the microscopic time scale (which is short compared to the macroscopic time).

We assume therefore that the molecular dynamics system can, locally in microscopic sets  $B$ , be viewed as a system in local equilibrium influenced by an external potential due to particle interaction outside the set, with a mean velocity and a temperature that can vary on a macroscopic space and time scale but are considered to be constant in the microscopic simulation set  $B = B_y$ , for microscopic time. To be in such local equilibrium is an approximation of a large system and the accuracy of this assumption depends on how fast  $T, u$  and  $R$  vary. The sets, that may overlap, move with the mean flow, following the *macroscopic characteristics* now defined by  $y^t = u(y^t, t)$ . The non interacting particle, whose position  $X_0$  measures time  $t$ , makes a well defined Hamiltonian system (11.14) also with given time dependent functions  $u = u(y^t, t)$  and  $R = R(y^t + X_j^t, t)$ .

The external field  $R$  is the potential due to particles outside the set  $B_y$  interacting with particles inside the set. The external potential  $R$  may depend on macroscopic time and we consider it as a boundary condition, acting on particles near the boundary, to get the right stress (and pressure) in a varying volume. The mean velocity implies that the particle positions  $X_j^t \in \mathbb{R}^3$  in Section ?? are replaced by  $X_j^t + y^t$ . Therefore, the external potential  $R(y^t + X_j^t, t) = \sum_{i \in B_y^c} \Phi(X_j^t + y^t - X_i - y^t)$  has the  $t$  dependence to capture the effect of the  $X_i^t + y^t$  dynamics. The molecular system (11.14) with an external force corresponds to simulation at given stress, instead of given volume  $\hat{V}$  in the standard  $NVT$  setting. In an equilibrium ensemble with constant pressure and varying volume, the quantity  $H_B = \hat{H} + \hat{P}\hat{V}$  (called the enthalpy in statistical mechanics and thermodynamics) replaces  $\hat{H}$ , where  $\hat{P}\hat{V} := -\sum_{j \in B} R(y^t + X_j^t, t)$  in a case when the

stress  $\sigma_{mn} = \hat{P}\delta_{mn}$  is given by the pressure.

We assume the local system in  $B_y$  is in *local equilibrium* - that is its probability density is the local Gibbs measure

$$G(X, p)dXd p := \frac{e^{-H_{B_y}(X, p)/T}}{\int e^{-H_{B_y}(X, p)/T} dXd p} dXd p. \quad (11.15)$$

We shall study the temperature  $T$ , velocity  $u$  and the density  $\rho : \mathbb{R}^3 \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$  as functions of the macroscopic space and length scales using molecular dynamics. We hence assume that electron-nuclei system locally relaxes to its local equilibrium (11.15), with macroscopic varying  $T$  and  $u$ . We will now study the effect of not being in macroscopic equilibrium, i.e. we ask what is the evolution of  $R$ ,  $T$  and  $u$ ? The answer will be the partial differential equation for compressible fluids describing the dynamics of the density  $\rho$ , the velocity  $u$  and the total energy  $E : \mathbb{R}^3 \times \mathbb{R}_+ \rightarrow \mathbb{R}$  in the system of the three *conservation laws*

$$\begin{aligned} \partial_t \rho + \partial_x(\rho u) &= 0 \\ \partial_t(\rho u) + \partial_x(\rho u \otimes u + \sigma) &= 0 \\ \partial_t E + \partial_x(Eu + \sigma \cdot u) &= 0 \end{aligned} \quad (11.16)$$

of mass, momentum and energy. To close the system one needs to relate the stress tensor  $\sigma : \mathbb{R}^3 \times \mathbb{R} \rightarrow \mathbb{R}^9$  with the other variables. The total energy can be written as a sum of kinetic energy and internal energy  $E = \rho|u|^2/2 + e$ , which defines the internal energy  $e$ . The stress tensor

$$\sigma = \sigma(e, \rho, u) \quad (11.17)$$

resulting from the "external" field  $R$  and the density, is a function of the internal energy  $e : \mathbb{R}^3 \times \mathbb{R}_+ \rightarrow \mathbb{R}$ , the density  $\rho$  and the velocity  $u$ ; the constitutive relation (11.17) can be determined from molecular dynamics simulations for a given fluid, e.g. using the microscopic formulation of internal energy and stress below. In the case of an ideal fluid  $\sigma_{nm} = c\rho e\delta_{nm}$ , for a constant  $c$ . The conservation of momentum can be written component wise as  $\partial_t(\rho u_i) + \sum_{j=1}^3 \partial_{x_j}(\rho u_i u_j + \sigma_{ij}) = 0$ . We use the notation  $\partial_x(ru) := \sum_{j=1}^3 \partial_{x_j}(ru_j)$  in the conservation of mass and energy equations.

Let  $\eta : \mathbb{R}^3 \rightarrow \mathbb{R}_+$  be a function which varies on the microscopic scale, has total integral  $\int_{\mathbb{R}^3} \eta(x) dx = 1$  and is supported on a tiny domain in the macroscopic scale, see Figure ??; hence  $\eta$  is an approximate delta-mass centered at the origin. The macroscopic *density*  $\rho : \mathbb{R}^3 \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$  of particles is defined by

$$\rho(x, t) := \int \sum_j \eta(x - X_j^t - y^t) G(X^0, p^0) dX^0 dp^0, \quad (11.18)$$

which we write as

$$\rho(x, t) = \mathbb{E}[\sum_j \eta(x - y^t - X^t)].$$

Smooth averages have been used in molecular dynamics for fluid dynamics, cf. [?], and for the vortex blob method and the smoothed particle hydrodynamics approximation of moving particles in fluid dynamics, cf. [?], [?]. We have

$$\partial_t \rho(x, t) = \mathbb{E} \left[ \sum_j \frac{d}{dt} \eta(x - X_j^t - y^t) \right]$$

and obtain by differentiation

$$\begin{aligned} \partial_t \rho(x, 0) &= - \int \sum_j \eta'(x - X_j - y) (\dot{X}_j + \dot{y}^0) G dX dp \\ &= - \int \sum_j \eta'(x - X_j - y^0) p_j G dX dp \\ &= - \int \sum_j \eta'(x - X_j - y^0) u(y^0, 0) G dX dp \end{aligned}$$

since

$$\int p_j G dp = u \int G dp = u \quad (11.19)$$

and  $\eta(x - X_j - y)$  is independent of  $p$ . By assumption the macroscopic function  $u$  is almost constant in the domain where  $\eta'$  is non zero, and we have obtained the conservation law of mass

$$\begin{aligned} \partial_t \rho(x, 0) &= -\partial_x \left( \int \sum_j \eta(x - X_j - y^0) u(x, 0) G dX dp \right) \\ &\quad + \partial_x \left( \int \sum_j \eta(x - X_j - y^0) (u(x, 0) - u(y^0, 0)) G dX dp \right) \\ &\rightarrow -\partial_x (\rho u(x, 0)) \end{aligned}$$

in the limit as  $\eta$  becomes a point mass in the macroscopic scale.

The next step is to derive the conservation law for momentum by differentiating the microscopic momentum

$$\int \sum_j \eta(x - X_j^t - y^t) p_j^t G(X^0, p^0) dX^0 dp^0 = \rho u(x, t).$$

We have similarly as for the density, using the special property of the Gibbs equilibrium density

$$T \partial_{X_j} G = -\partial_{X_j} H_B G = - \underbrace{(\partial_{X_j} \hat{H} - R'(y + X_j))}_{=-\dot{p}_j} G, \quad (11.20)$$

that

$$\begin{aligned}
\partial_t(\rho u)(x, 0) &= - \int \sum_j (\eta'(x - X_j - y^0) p_j (\dot{X}_j + \dot{y}) - \eta(x - X_j - y^0) \dot{p}_j) G dX dp \\
&= - \int \sum_j (\eta'(x - X_j - y^0) p_j \otimes p_j G - \eta(x - X_j - y^0) \dot{p}_j) G dX dp \\
&= - \int \sum_j \left( \eta'(x - X_j - y^0) p_j \otimes p_j G \right. \\
&\quad \left. - \eta(x - X_j - y^0) T \partial_{X_j} G + \eta(x - X_j - y_0) \partial_{X_j} R(y^0 + X_j, 0) G \right) dX dp.
\end{aligned}$$

The integration by parts, using (11.20), is called a *virial property*. By writing  $p_j = (p_j - u) + u$  and using  $\int (p_j - u) G dp = 0$  together with

$$\int (p_j^n - u_n)(p_j^m - u_m) G dp = \begin{cases} T & n = m \\ 0 & n \neq m \end{cases}, \quad (11.21)$$

we have

$$\begin{aligned}
\partial_t(\rho u) &= - \int \sum_j \left( \eta'(x - X_j - y)(u \otimes u + T) \right. \\
&\quad \left. + \partial_{X_j}(\eta(x - X_j - y)T) + \eta(x - X_j - y) \partial_{X_j} R(y + X_j, 0) \right) G dX dp \\
&= - \partial_x \int \sum_j \left( \eta(x - X_j - y)(u \otimes u + T) - \eta(x - X_j - y)T \right) G dX dp \quad (11.22) \\
&\quad - \int \sum_j \eta(x - X_j - y) \partial_y R(y + X_j, 0) G dX dp.
\end{aligned}$$

We want a spacial derivative with respect to  $x$  on the last integral of the external forces to get the conservative stress field. This can be obtained from the construction

$$\zeta(x - y, X_j, X_i) := \int_0^1 \eta(x - y - X_j + \lambda(X_j - X_i)) d\lambda,$$

since we have

$$\eta(x - y - X_j) - \eta(x - y - X_i) = (X_j - X_i) \partial_x \zeta(x - y, X_j, X_i)$$

and for particles  $i$  in  $B_y^c$  there holds  $\eta(x - y - X_i) = 0$ , so that

$$\begin{aligned}
& \int \sum_j \eta(x - X_j - y) \partial_y R(y + X_j, 0) G dX dp \\
&= \int \sum_j \sum_{i \in B_y^c} \eta(x - X_j - y) \Phi'(X_j - X_i) G dX dp \\
&= \int \sum_j \sum_{i \in B_y^c} (\underbrace{\eta(x - y - X_j) - \eta(x - y - X_i)}_{=0}) \Phi'(X_j - X_i) G dX dp \quad (11.23) \\
&= \partial_x \int \sum_j \sum_{i \in B_y^c} (X_j - X_i) \zeta(x - y, X_j, X_i) \Phi'(X_j - X_i) G dX dp \\
&=: \partial_x \sigma
\end{aligned}$$

defines the stress tensor  $\sigma : \mathbb{R}^3 \times \mathbb{R}_+ \rightarrow \mathbb{R}^9$ . We have obtained the conservation law of momentum

$$\partial_t(\rho u) + \partial_x(\rho u \otimes u + \sigma) = 0.$$

Note that the two pressure like terms  $\rho T$  and  $-\rho T$  in (11.22), from the fluctuation of the kinetic energy respectively from the interaction of particle forces, cancel each other. We see that the stress tensor is symmetric for a potential depending on the pair distance, since

$$\Phi'(X_j - X_i) = \partial_{X_j} \tilde{\Phi}(|X_j - X_i|) = (X_j - X_i) \tilde{\Phi}'(|X_j - X_i|) / |X_j - X_i|.$$

As usual the pressure  $P$  is one third of the trace of the stress tensor.

The final step is to derive the conservation of energy by differentiation of the microscopic *total energy*

$$E := \int \sum_{j \in B_y} \eta(x - X_j - y) \left( \frac{|p_j|^2}{2} + \frac{1}{2} \sum_{i \in B_y, i \neq j} \Phi(X_j - X_i) \right) G(X, p) dX dp, \quad (11.24)$$

which has the *kinetic energy* part

$$\int \sum_j \eta(x - X_j - y) \frac{|p_j|^2}{2} G dX dp = \frac{\rho |u|^2}{2} + \frac{3\rho T}{2}$$

and the *potential energy* part

$$\frac{1}{2} \int \left( \sum_j \eta(x - X_j - y) \sum_{i \in B_y, i \neq j} \Phi(X_j - X_i) \right) G(X, p, t) dX dp =: m. \quad (11.25)$$

The reason we use a pair potential is that it allows for the simple interpretation of the potential energy related to one particle presented in (11.25); Section 11.4.1 describes a

generalization from the pair potential  $2^{-1} \sum_j \sum_{i \neq j} \Phi(X_j - X_i)$  to an arbitrary potential  $U(X)$ . Let the internal energy be the sum of these kinetic and potential energies

$$e := \frac{3\rho T}{2} + m. \quad (11.26)$$

We have as above

$$\begin{aligned} & \partial_t \left( \frac{\rho |u|^2}{2} + e \right) \\ &= - \int \sum_j (\eta'(x - X_j - y) \left( \frac{|p_j|^2}{2} + \frac{1}{2} \sum_{i \neq j} \Phi(X_j - X_i) \right) p_j G dX dp \\ & \quad + \int \sum_j \eta(x - X_j - y) \left( \dot{p}_j \cdot p_j + \frac{1}{2} \sum_{i \neq j} \Phi'(X_j - X_i) (p_j - p_i) \right) G dX dp \\ & \rightarrow -\partial_x \left( \rho \left( \frac{|u|^2}{2} + \frac{3T}{2} \right) u + mu + \sigma \cdot u \right) \\ &= -\partial_x \left( \left( \rho \left( \frac{|u|^2}{2} + \frac{3T}{2} \right) + m \right) u + \sigma \cdot u \right) \end{aligned}$$

using  $\int (p_j - p_i) G dp = 0$  and  $p = p - u + u$  (to the second and third power) to obtain  $\int |p|^2 p G dp = (|u|^2 + 3T + 2T)u$ . We conclude that the total energy  $E := \rho |u|^2/2 + e$  satisfies the conservation law

$$\partial_t E + \partial_x (Eu + \sigma \cdot u) = 0.$$

The derivation of the macroscopic equations (11.16) of compressible flow also gave us microscopic definitions of the bulk density (11.18), velocity (11.19), temperature (11.21), stress (11.23), energy (11.24) and internal energy (11.25-11.26).

### 11.4.1 A general potential

Consider a case with a general potential  $U(X)$  replacing  $2^{-1} \sum_j \sum_{i \neq j} \Phi(X_j - X_i)$  in (11.14). The derivation of the conservation laws for momentum uses that the stress can be defined from the interaction with particles outside the set, which is simple for a pair potential. For a general potential one need also to identify such interactions outside the set. Therefore we split the potential. The splitting assumes that  $U$  can be split into potential energies related to the individual particles

$$U(X) = \sum_j m_j(X),$$

where each term  $m_j$  corresponds to  $\sum_{i \neq j} \Phi(X_j - X_i)/2$  in the pair potential case. To split the Gibbs measure into local equilibrium parts, let

$$U_B := \sum_{j \in B} m_j(X),$$



which defines the external part

$$U - U_B =: R_1$$

and yields

$$\partial_{X_j} U = \partial_{X_j} (U_B + R_1).$$

Assume we can split the external part into contributions from particles outside  $B$

$$R_1 = \sum_{i \in B^c} r_i.$$

The stress is then defined with  $\partial_{X_j} r_i(X)$  replacing  $\Phi'(X_j - X_i)$  in (11.23) and the local equilibrium measure

$$\frac{e^{-(\sum_{j \in B} |p_j|^2/2 + U_B(X))/T} dX dp}{\int e^{-(\sum_{j \in B} |p_j|^2/2 + U_B(X))/T} dX dp}$$

replaces  $e^{-H_B/T} dX dp / \int e^{-H_B/T} dX dp$  in (11.15).

To handle the conservation laws of energy, the derivation with pair potentials uses in addition in fact only that

$$\frac{d}{dt} \int \frac{1}{2} \sum_{i \neq j} \Phi(X_j^t - X_i^t) G(X^0, p^0) dp^0 = \int \frac{1}{2} \sum_{i \neq j} \Phi'(X_j^0 - X_i^0) \cdot (p_j^0 - p_i^0) G dp^0 = 0$$

which follows from  $\int (p_j - p_i) G dp = 0$ . Suppose we have a partition  $U(X) = \sum_j m_j(X)$  that satisfies

$$\sum_k \int u \cdot \partial_{X_k} m_j(X) G dX = 0$$

for every  $j$ . Then the derivation above can be applied, with  $m_j(X)$  replacing  $\sum_{i \neq j} \Phi(X_j - X_i)$ .

With a derivation based on a general molecular dynamics potential we can for instance use the Ehrenfest dynamics with the Hamiltonian  $|p - u|^2/2 + \phi \cdot V(X)\phi$  and the corresponding equilibrium measure

$$G dX dp d\phi_r d\phi_i = e^{-(|p-u|^2/2 + \phi \cdot V'(X)\phi)/T} dX dp d\phi_r d\phi_i.$$

Using that the Ehrenfest dynamics is a Hamiltonian system in the variables  $X, p$  and  $\phi = \phi_r + i\phi_i$ , also the virial term can be handled as above since the  $\phi \cdot V'(X)\phi = -T \partial_X G$ .

## Chapter 12

# Appendices

### 12.1 Tomography Exercise

Tomographic imaging is used in medicine to determine the shape/image of a bone or interior organ. One procedure for doing this is by projecting X-rays from many different angles through the body (see figure 1), measure the strength of the X-rays that has gone through the image, and compute how the image has to be to comply with the X-ray output. Reconstructing an image this way is called tomographic reconstruction, and it is the problem we look at in this project.

In our case we first superimpose a grid over the image we wish to perform tomographic imaging on to an  $n \times n$  pixel image represented with image values as vector  $(f_i)_{i=1}^{n^2}$ . The image values are assumed to be constant within each cell of the grid. An  $n = 3$  case with vertical and horizontal projections serves the purpose of further explaining the problem: In figure 2 we have superimposed a  $3 \times 3$  square grid on the image  $f(x, y)$ . The rays are the lines running through the  $x - y$  plane (we disregard the width of the lines here assuming they are all of the same width and very thin). The projections are given the representation  $p_i$ , we say that  $p_i$  is the ray sum measured with the  $i$ th ray. The relationship between the  $f_j$ 's and the  $p_i$ 's may be expressed as the set of linear equations

$$\sum_{j=1}^{n^2} A_{ij} f_j = p_i, \quad i = 1, \dots, n. \quad (12.1)$$

For example, the first equation in the  $3 \times 3$  case only goes through  $f_1, f_4$  and  $f_7$  yielding the equation

$$A_{11} f_1 + A_{14} f_4 + A_{17} f_7 = p_1,$$

The linear system of equations created by the horizontal and vertical projections in figure

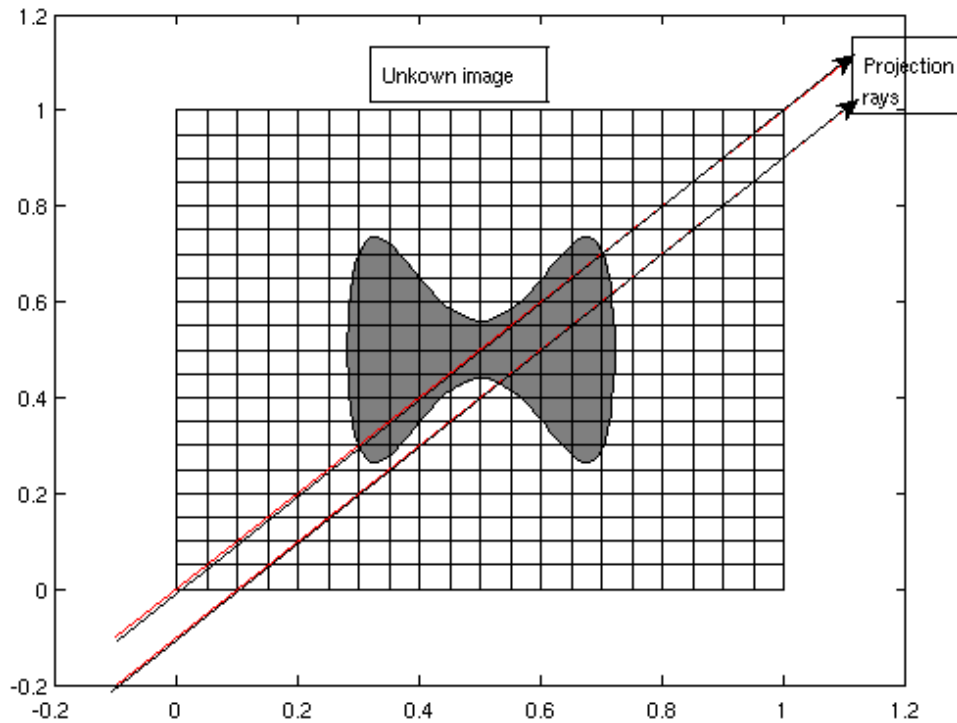


Figure 12.1: Illustration of tomographic imaging. The image on the unit square represents our unknown image which we send rays through to determine.

2 written on the form  $An = p$  is

$$\begin{pmatrix} 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \\ f_5 \\ f_6 \\ f_7 \\ f_8 \\ f_9 \end{pmatrix} = \begin{pmatrix} p_1 \\ p_2 \\ p_3 \\ p_4 \\ p_5 \\ p_6 \end{pmatrix} \quad (12.2)$$

In this case,  $A \in \mathbb{R}^{6 \times 9}$ . The problem is underdetermined so the least squares way of solving this problem:

$$f = (A^T A)^{-1} A^T p, \quad (12.3)$$

fails because  $A^T A$  is singular. One way to deal with the singular matrix is to instead

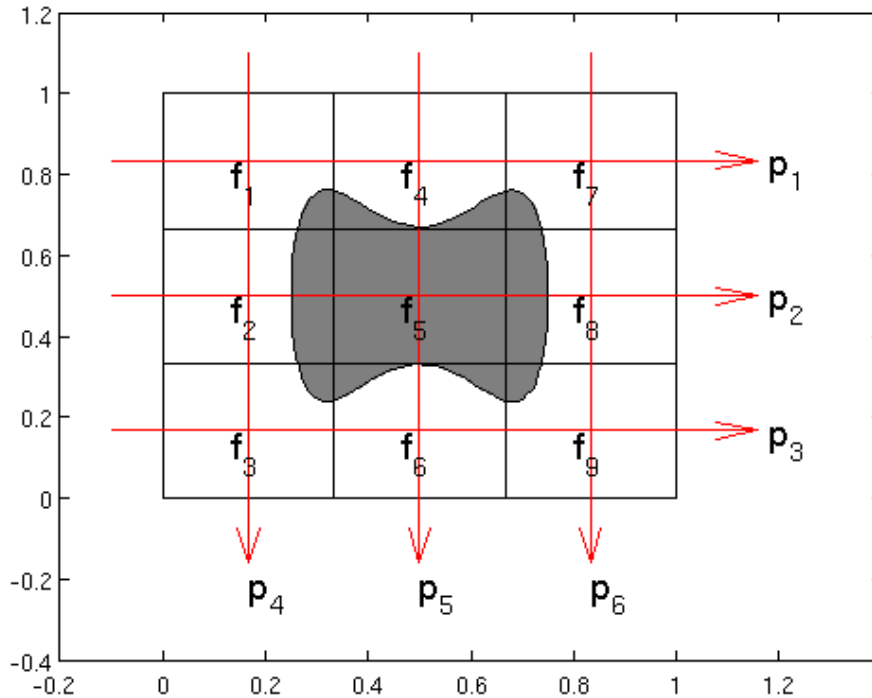


Figure 12.2: Illustration of horizontal and vertical projections on a  $3 \times 3$  image.

solve

$$f = (A^T A + \delta I_{n^2})^{-1} A^T p,$$

where  $\delta$  is a small number.

### Exercise 1.

Download the image “ImageEx1.jpg” and the matlab program “rayItHorVert.m”. This image is our unknown image (we only have the solution to compare). Create an image matrix by the command

```
image = imread('ImageEx1.tif')
```

Create a projection vector of the image by calling

```
p=rayItHorVert(f)
```

Write a matlab program that takes as input a vector  $p \in \mathbb{R}^{6 \times 1}$ , creates the matrix  $A \in \mathbb{R}^{6 \times 9}$  given in (12.2) (for  $n = 3$ ) and finds the tomographically reconstructed image  $f$  by the computation (12.3). Use

```
f=reshape(f,n,n)
```

to reshape the vector  $f$  into an  $n \times n$  matrix and plot by the commands

```
colormap(gray)
imagesc(f)
```

Also plot the matrix “image” and compare results. As a reference, the result should look like figure 3:

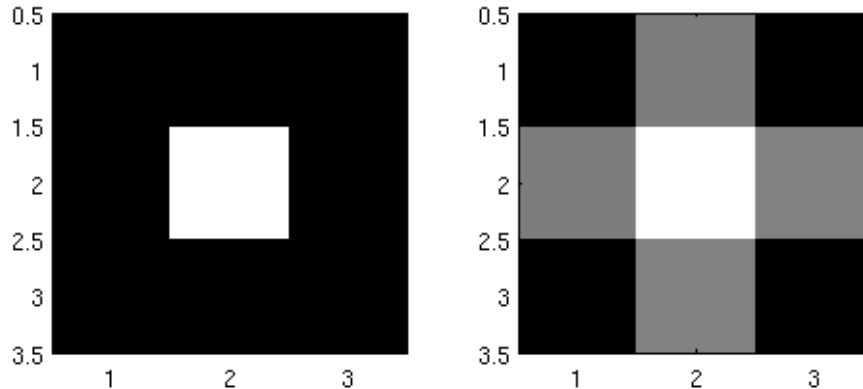


Figure 12.3: Illustration of the image “ImageEx1.jpg” (left) and the tomographic reconstruction (right).

*Hint: The matrix  $A$  can be created quite easily with the Kronecker product  $\otimes$  which is defined as follows:*

$$B \otimes C = \begin{pmatrix} BC_{11} & BC_{12} & \dots & BC_{1n} \\ BC_{21} & BC_{22} & \dots & BC_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ BC_{m1} & BC_{m2} & \dots & BC_{mn} \end{pmatrix} \quad (12.4)$$

where  $C \in \mathbb{R}^{m \times n}$  and  $B$  is an arbitrary matrix. In matlab the operation  $B \otimes C$  is written `kron(B,C)`

### Exercise 2.

Use the hint in exercise 1. to generalize the matlab program to work for any  $n$  value. That is, write a program that takes as input an  $n$ -value and a vector  $p \in \mathbb{R}^{2n \times 1}$ , and creates a matrix  $A \in \mathbb{R}^{2n \times n^2}$  with similar structure as the one in (12.2).

(a)

Download the image “Ball.tif” and solve the problem as in exercise one. One might improve the reconstructed image quality by filtering the image. Implement a scheme which removes values below a certain threshold in the matrix  $f$  and plot the result.

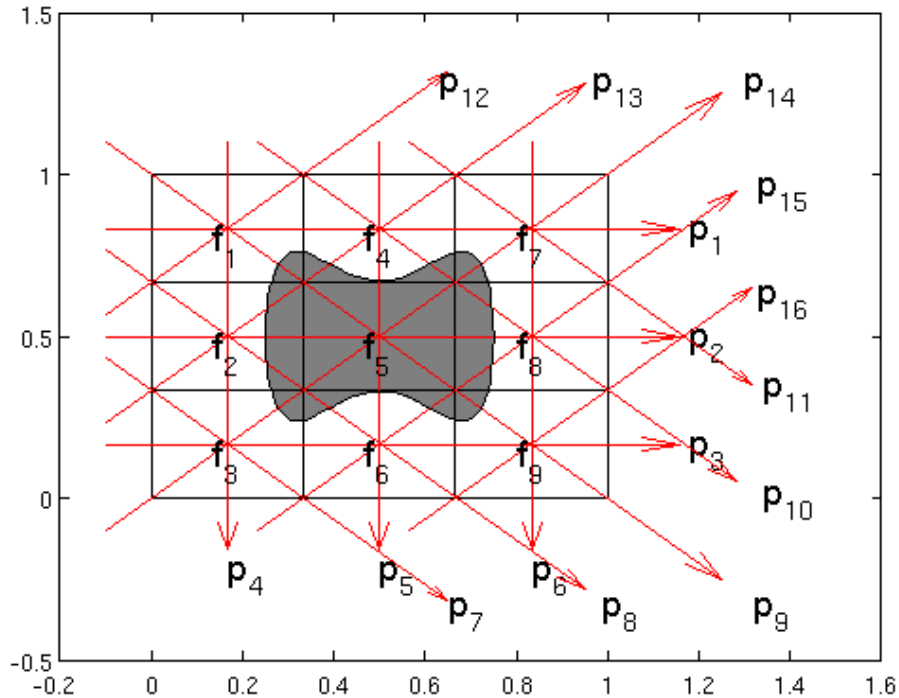


Figure 12.4: Illustration of horizontal, vertical and diagonal projections on a  $3 \times 3$  image.

(b)

Assume that you have X-rayed a square shaped suitcase containing a circular shaped bomb. The file “pVector.mat” consists of the projection vector which you read by the command

```
load('pVector.mat')
```

What is approximately the position of the bomb? (Assume unit square coordinates).

(c)

Download the image “TwoBalls.tif” and solve the problem as in exercise one. Why does the reconstructed image differ so strongly from the real one?

The scheme implemented in exercise 3 improves the reconstructed image.

### Exercise 3. - Week project exercise

The next step is to add more projections to our tomographic imaging. As illustrated in figure 4, we use horizontal, vertical and diagonal projections. For the  $n = 3$  case the

linear set of equations  $Af = p$  is

$$\begin{pmatrix}
 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\
 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\
 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 \\
 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\
 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\
 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\
 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1
 \end{pmatrix}
 \begin{pmatrix}
 f_1 \\
 f_2 \\
 f_3 \\
 f_4 \\
 f_5 \\
 f_6 \\
 f_7 \\
 f_8 \\
 f_9
 \end{pmatrix}
 =
 \begin{pmatrix}
 p_1 \\
 p_2 \\
 p_3 \\
 p_4 \\
 p_5 \\
 p_6 \\
 p_7 \\
 p_8 \\
 p_9 \\
 p_{10} \\
 p_{11} \\
 p_{12} \\
 p_{13} \\
 p_{14} \\
 p_{15} \\
 p_{16}
 \end{pmatrix}
 \tag{12.5}$$

Write a program that takes as input an  $n$ -value and a vector  $p \in \mathbb{R}^{(6n-2) \times 1}$ , and creates a matrix  $A \in \mathbb{R}^{(6n-2) \times n^2}$  with similar structure to the one in (12.5). Download the image “TwoBalls.tif” and the program “rayItHorVertDiag.m” which you use to create the projection vector by the command

```
p=rayItHorVertDiag(f)
```

Solve this image problem as in exercise 2 (c). Implement the filtering technique here as well. Compare this reconstruction to the one in 2 (c).

#### Exercise 4. - Week project exercise

The reason we are looking at low resolution images above is that for an  $n \times n$  image the matrix  $A \in \mathbb{R}^{(6n-2) \times n^2}$ . This means that  $A^T A \in \mathbb{R}^{n^2 \times n^2}$  which is so huge, that even for relatively small  $n$  that we can not solve the problem (12.3) in Matlab the way we have done in the exercises above. The paper “Algebraic reconstruction algorithms” describes an iterative algorithm solving the tomographic reconstruction problem which works for higher resolution images (see page 278). Read the first pages of this paper and implement this algorithm using horizontal, vertical and diagonal projections as in exercise 3. Try your algorithm on the picture “Pear.tif”

## 12.2 Molecular Dynamics

Here some discussion about the MD code will appear.

```
#include <math.h>
#include <stdio.h>
```

```

#include <stdlib.h>
#include <iostream>
#include <iomanip>
#include <fstream>
#include <sstream>
#include <string>

//
// Compile with g++ -O2 -o main main.cpp
//

using namespace std;

// ----- Definitions -----
typedef double real;
real sqr(real n){return n*n;}
enum BoundaryCond {periodic, flow};

// ----- Cell and particle structures -----
struct Parameters
{
    real sigma, epsilon, cutoff, dt, T, temp, size[3];
    int cells[3], cellsTot;
    BoundaryCond bc;
};

struct Particle
{
    real m;
    real x[3];
    real v[3];
    real F[3];
    real Fold[3];
    int flag;
};

struct ParticleList
{
    Particle p;
    ParticleList *next;
};

typedef ParticleList* Cell;

```



```

void insertListElement(ParticleList **root, ParticleList *pl)
{
    pl->next = *root;
    *root = pl;
}

void deleteListElement(ParticleList **pl)
{
    *pl = (*pl)->next;
}

int index(int *i, int *cells)
{
    return i[0] + cells[0]*(i[1] + cells[1]*i[2]);
}

// ----- Function definitions -----
void inputParameters(Parameters&);
void initData(Cell*, Parameters&);
void integrate(real, Cell*, Parameters&);
void compF(Cell*, Parameters&);
void compX(Cell*, Parameters&);
void compV(Cell*, Parameters&);
real compE(Cell*, Parameters&);
void updateX(Particle*, real);
void updateV(Particle*, real);
void forceLJ(Particle*, Particle*, real, real);
void sortParticles(Cell*, Parameters&);
void saveParticles(Cell*, real, Parameters&);
void boltzmann(Particle*, real);
real gaussDeviate();

// ----- Program and functions -----
int main(int argc, char **argv)
{
    int s = system("rm -rf ./data/*.txt");
    Parameters p;
    inputParameters(p);
    Cell *grid = new Cell[p.cellsTot];
    //for (int i=0; i<p.cellsTot; ++i)
    // grid[i] = NULL;
    initData(grid, p);
}

```

```

    saveParticles(grid, 0, p);
    integrate(0, grid, p);
    return s;
}

void inputParameters(Parameters &p)
{
    // Lennard-Jones parameters
    p.sigma = 3.4;
    p.epsilon = 120;
    // Box size
    for (int d=0; d<3; ++d)
        p.size[d] = 150*p.sigma;
    // Cells
    p.cutoff = 2.5*p.sigma;
    for (int d=0; d<3; ++d)
        p.cells[d] = (int) floor(p.size[d] / p.cutoff);
    p.cellsTot = 1;
    for (int d=0; d<3; ++d)
        p.cellsTot *= p.cells[d];
    // Timescale
    p.T = 20;
    p.dt = 1e-2;
    // Boundary condition
    p.bc = flow;
    // Save to file
    FILE *file = fopen("./data/parameters.txt", "w");
    fprintf(file, "%f %f %f %f %f ", p.sigma, p.epsilon, p.cutoff, p.T, p.dt);
    for (int d=0; d<3; ++d)
        fprintf(file, "%f ", p.size[d]);
    for (int d=0; d<3; ++d)
        fprintf(file, "%d ", p.cells[d]);
    fclose(file);
}

void initData(Cell *grid, Parameters &p)
{
    // Box 1
    real mass = 39.95;
    int n1 = 10, n2 = 10, n3 = 10;
    grid[0] = NULL;
    ParticleList **root = &grid[0];
}

```

```

for (int i=0; i<=2*n1; ++i)
  for (int j=0; j<=2*n2; ++j)
    for (int k=0; k<=2*n3; ++k)
      {
        // Face centered cubic
        if ( !((i+j+k)%2) )
          {
            ParticleList *pl = new ParticleList;
            pl->p.m = mass;
            pl->p.x[0] = 0.5*p.size[0] + (i-n1)*pow(2, 1.0/6.0)*p.sigma;
            pl->p.x[1] = 0.5*p.size[1] + (j-n2)*pow(2, 1.0/6.0)*p.sigma;
            pl->p.x[2] = 0.6*p.size[2] + (k-n3)*pow(2, 1.0/6.0)*p.sigma;
            pl->p.v[0] = 0;
            pl->p.v[1] = 0;
            pl->p.v[2] = -20.4;
            pl->p.flag = 0;
            insertListElement(root, pl);
          }
      }

// Box 2
n1 = 30, n2 = 30, n3 = 10;
for (int i=0; i<=2*n1; ++i)
  for (int j=0; j<=2*n2; ++j)
    for (int k=0; k<=2*n3; ++k)
      {
        // Face centered cubic
        if ( !((i+j+k)%2) )
          {
            ParticleList *pl = new ParticleList;
            pl->p.m = mass;
            pl->p.x[0] = 0.5*p.size[0] + (i-n1)*pow(2, 1.0/6.0)*p.sigma;
            pl->p.x[1] = 0.5*p.size[1] + (j-n2)*pow(2, 1.0/6.0)*p.sigma;
            pl->p.x[2] = 0.4*p.size[2] + (k-n3)*pow(2, 1.0/6.0)*p.sigma;
            pl->p.v[0] = 0;
            pl->p.v[1] = 0;
            pl->p.v[2] = 0;
            pl->p.flag = 1;
            insertListElement(root, pl);
          }
      }

// Noise
for (ParticleList *pl=grid[0]; pl!=NULL; pl=pl->next)

```

```

    boltzmann(&p1->p, 1.0);
    sortParticles(grid, p);
}

void boltzmann(Particle *p, real factor)
{
    for (int d=0; d<3; ++d)
        p->v[d] += factor * gaussDeviante();
}

real gaussDeviante()
{
    real a1, a2, s, r, b1;
    static int iset = 0;
    static real b2;
    if (!iset)
    {
        do {
            a1 = 2.0 * rand() / (RAND_MAX + 1.0) - 1.0;
            a2 = 2.0 * rand() / (RAND_MAX + 1.0) - 1.0;
            r = a1 * a1 + a2 * a2;
        } while (r>=1.0);
        s = sqrt(-2.0 * log(r) / r);
        b1 = a1 * s;
        b2 = a2 * s;
        iset = 1;
        return b1;
    }
    else
    {
        iset = 0;
        return b2;
    }
}

void integrate(real t, Cell *grid, Parameters &p)
{
    compF(grid, p);
    while (t < p.T-1e-9)
    {
        t += p.dt;
        compX(grid, p);
        compF(grid, p);
        compV(grid, p);
        saveParticles(grid, t, p);
    }
}

```

```

        cout << scientific <<
            "t = " << t << " E = " << compE(grid, p) << endl;
    }
}

void compF(Cell *grid, Parameters &p)
{
    int* cells = p.cells;
    int i[3], j[3];
    // Loop over cells in each dimension
    for (i[0]=0; i[0]<cells[0]; i[0]++)
        for (i[1]=0; i[1]<cells[1]; i[1]++)
            for (i[2]=0; i[2]<cells[2]; i[2]++)
                // Loop over particles in each cell
                for (ParticleList *p11=grid[index(i,cells)]; p11!=NULL; p11=p11->next)
                    {
                        for (int d=0; d<3; ++d)
                            p11->p.F[d] = 0;
                        // Loop over neighbours in each dimension
                        for (j[0]=i[0]-1; j[0]<=i[0]+1; j[0]++)
                            for (j[1]=i[1]-1; j[1]<=i[1]+1; j[1]++)
                                for (j[2]=i[2]-1; j[2]<=i[2]+1; j[2]++)
                                    {
                                        bool outside = false;
                                        int tmp[3];
                                        if (p.bc==periodic)
                                            {
                                                // Periodic boundary
                                                for (int d=0; d<3; ++d)
                                                    tmp[d] = j[d];
                                                for (int d=0; d<3; ++d)
                                                    if (j[d]<0)
                                                        j[d] = cells[d]-1;
                                                    else if (j[d]>=cells[d])
                                                        j[d] = 0;
                                            }
                                        else if (p.bc==flow)
                                            {
                                                // Flow boundary
                                                for (int d=0; d<3; ++d)
                                                    if (j[d]<0 || j[d]>=cells[d])
                                                        outside = true;
                                            }
                                        if (!outside)

```

```

        {
            // Check distance from particle pl1 to neighbour cell j
            real dist = 0;
            for (int d=0; d<3; ++d)
                dist +=
                    sqr( min( pl1->p.x[d] - j[d] * 1.0 / cells[d],
                             pl1->p.x[d] - (j[d] + 1) * 1.0 / cells[d] ) );
            // Loop over particles in each neighbour cell
            //if (dist<=p.cutoff)
            for (ParticleList *pl2=grid[index(j,cells)]; pl2!=NULL; pl2=pl2->next)
                if (pl1!=pl2)
                    {
                        real r = 0;
                        for (int d=0; d<3; ++d)
                            r += sqr(pl2->p.x[d] - pl1->p.x[d]);
                        if (r<=sqr(p.cutoff))
                            forceLJ(&pl1->p, &pl2->p, p.sigma, p.epsilon);
                    }
        }
        if (p.bc==periodic)
            {
                // Copy back
                for (int d=0; d<3; ++d)
                    j[d] = tmp[d];
            }
    }
}

void forceLJ(Particle *i, Particle *j, real sigma, real epsilon)
{
    real r = 0.0;
    for (int d=0; d<3; ++d)
        r += sqr(j->x[d] - i->x[d]);
    real s = sqr(sigma) / r;
    s = sqr(s) * s;
    real f = 24 * epsilon * s / r * (1 - 2 * s);
    for (int d=0; d<3; ++d)
        i->F[d] += f * (j->x[d] - i->x[d]);
}

void compX(Cell *grid, Parameters &p)
{
    int i[3];

```

```

// Loop over cells in each dimension
for (i[0]=0; i[0]<p.cells[0]; i[0]++)
  for (i[1]=0; i[1]<p.cells[1]; i[1]++)
    for (i[2]=0; i[2]<p.cells[2]; i[2]++)
      // Loop over particles in each cell
      for (ParticleList *pl=grid[index(i,p.cells)]; pl!=NULL; pl=pl->next)
        updateX(&pl->p, p.dt);

// Update cells according to new positions
sortParticles(grid, p);
}

void updateX(Particle *p, real dt)
{
  real a = dt * 0.5 / p->m;
  for (int d=0; d<3; ++d)
    {
      p->x[d] += dt * (p->v[d] + a * p->F[d]);
      p->Fold[d] = p->F[d];
    }
}

void compV(Cell *grid, Parameters &p)
{
  int i[3];
  // Loop over cells in each dimension
  for (i[0]=0; i[0]<p.cells[0]; i[0]++)
    for (i[1]=0; i[1]<p.cells[1]; i[1]++)
      for (i[2]=0; i[2]<p.cells[2]; i[2]++)
        // Loop over particles in each cell
        for (ParticleList *pl=grid[index(i,p.cells)]; pl!=NULL; pl=pl->next)
          updateV(&pl->p, p.dt);
}

void updateV(Particle *p, real dt)
{
  real a = dt * 0.5 / p->m;
  for (int d=0; d<3; ++d)
    {
      p->v[d] += a * (p->F[d] + p->Fold[d]);
    }
}

void sortParticles(Cell *grid, Parameters &p)

```

```

{
  int i[3], j[3];
  // Loop over cells in each dimension
  for (i[0]=0; i[0]<p.cells[0]; i[0]++)
    for (i[1]=0; i[1]<p.cells[1]; i[1]++)
      for (i[2]=0; i[2]<p.cells[2]; i[2]++)
        {
          // Pointers to particle list in cell i
          ParticleList **pl1 = &grid[index(i,p.cells)];
          ParticleList *pl2 = *pl1;
          // Traverse list in cell i
          while (pl2!=NULL)
            {
              bool outside = false;
              // Cell that particle belongs to
              for (int d=0; d<3; ++d)
                {
                  j[d] = (int) floor(pl2->p.x[d] * p.cells[d] / p.size[d]);
                  if (p.bc==periodic)
                    {
                      // Periodic boundary
                      if (j[d]<0)
                        j[d] = p.cells[d] - j[d] % p.cells[d];
                      else if (j[d]>=p.cells[d])
                        j[d] = j[d] % p.cells[d];
                    }
                  else if (p.bc==flow)
                    {
                      // Outflow boundary
                      if (j[d]<0 || j[d]>=p.cells[d])
                        outside = true;
                    }
                }
              // If not same cell
              if ( (i[0]!=j[0]) || (i[1]!=j[1])
                  || (i[2]!=j[2]) )
                {
                  // Delete particle from list
                  deleteListElement(pl1);
                  // Add to list in cell j
                  if (!outside)
                    insertListElement(&grid[index(j,p.cells)], pl2);
                }
            }
          else

```



```

        pl1 = &pl2->next;
        pl2 = *pl1;
    }
}

real compE(Cell* grid, Parameters &p)
{
    real e = 0;
    int i[3];
    // Loop over cells in each dimension
    for (i[0]=0; i[0]<p.cells[0]; i[0]++)
        for (i[1]=0; i[1]<p.cells[1]; i[1]++)
            for (i[2]=0; i[2]<p.cells[2]; i[2]++)
                // Loop over particles in each cell
                for (ParticleList *pl=grid[index(i,p.cells)]; pl!=NULL; pl=pl->next)
                    {
                        real v = 0;
                        for (int d=0; d<3; ++d)
                            v += sqr(pl->p.v[d]);
                        e += 0.5 * pl->p.m * v;
                    }
    return e;
}

void saveParticles(Cell* grid, real t, Parameters &p)
{
    stringstream ss;
    ss.str(""); ss << fixed << setprecision(6) << t/p.T;
    string fname("./data/" + ss.str() + ".txt");
    FILE *file = fopen(fname.c_str(), "w");

    int i[3];
    // Loop over cells in each dimension
    for (i[0]=0; i[0]<p.cells[0]; i[0]++)
        for (i[1]=0; i[1]<p.cells[1]; i[1]++)
            for (i[2]=0; i[2]<p.cells[2]; i[2]++)
                {
                    // Loop over particles in each cell
                    for (ParticleList *pl=grid[index(i,p.cells)]; pl!=NULL; pl=pl->next)
                        {
                            for (int d=0; d<3; ++d)
                                fprintf(file, "%f ", pl->p.x[d]);
                            for (int d=0; d<3; ++d)

```

```
        fprintf(file, "%f ", pl->p.v[d]);
        fprintf(file, "%d \n", pl->p.flag);
    }
}
fclose(file);
}
```

# Chapter 13

## Recommended Reading

The following references have been useful for preparing these notes and are recommended for further studies.

### Stochastic Differential Equations

- Online material: [\[Evab\]](#)
- Numerics for SDE: [\[KP92, Mil95\]](#)
- SDE: [\[Øks98\]](#)
- Advanced SDE: [\[KS91\]](#)

### Probability

[\[Dur96\]](#)

### Mathematical Finance

- Basic stochastics for finance: [\[Bax96\]](#)
- Finance in practice: [\[Hul97\]](#)
- Finance with numerics: [\[WD95\]](#)

### Partial Differential Equations

- Advanced PDE: [\[Eva98\]](#)
- Online introduction: [\[Evaa\]](#)
- FEM: [\[Joh87\]](#)
- Advanced FEM: [\[BS94\]](#)
- Introductory DE and PDE: [\[EEHJ96\]](#) and [\[Str86\]](#)

## Variance Reduction for Monte Carlo Methods

[Caf98]

## Molecular Dynamics

[LB05], [CDK<sup>+</sup>03b], [Fre02]

# Bibliography

- [AVE09] G Ariel and E Vanden-Eijnden. A strong limit theorem in the  $\epsilon$ -oscillator model. *Nonlinearity*, 22(1):145–162, 2009.
- [Bax96] Andrew Baxter, Martinand Rennie. *Financial calculus : an introduction to derivative pricing*. Cambridge Univ. Press, Cambridge, 1996.
- [BBK07] John S. Briggs, Sutee Boonchui, and Supitch Khemmani. The derivation of time-dependent Schrödinger equations. *J. Phys. A*, 40(6):1289–1302, 2007.
- [BNS96] Folkmar A. Bornemann, Peter Nettesheim, and Christof Schütte. Quantum-classical molecular dynamics as an approximation to full quantum dynamics. *The Journal of Chemical Physics*, 105(3):1074–1083, 1996.
- [BO27] M. Born and R. Oppenheimer. Zur Quantentheorie der Molekeln. *Annalen der Physik*, 389:457–484, 1927.
- [BR01] John S. Briggs and Jan M. Rost. On the derivation of the time-dependent equation of Schrödinger. *Found. Phys.*, 31(4):693–712, 2001. Invited papers dedicated to Martin C. Gutzwiller, Part V.
- [BS91] F. A. Berezin and M. A. Shubin. *The Schrödinger equation*, volume 66 of *Mathematics and its Applications (Soviet Series)*. Kluwer Academic Publishers Group, Dordrecht, 1991. Translated from the 1983 Russian edition by Yu. Rajabov, D. A. Leites and N. A. Sakharova and revised by Shubin, With contributions by G. L. Litvinov and Leites.
- [BS94] Susanne C. Brenner and L. Ridgway Scott. *The mathematical theory of finite element methods*, volume 15 of *Texts in Applied Mathematics*. Springer-Verlag, New York, 1994.
- [Caf98] Russel E. Caflisch. Monte Carlo and quasi-Monte Carlo methods. In *Acta numerica, 1998*, volume 7 of *Acta Numer.*, pages 1–49. Cambridge Univ. Press, Cambridge, 1998.
- [CDK<sup>+</sup>03a] Eric Cancès, Mireille Defranceschi, Werner Kutzelnigg, Claude Le Bris, and Yvon Maday. Computational quantum chemistry: a primer. In *Handbook of numerical analysis, Vol. X*, Handb. Numer. Anal., X, pages 3–270. North-Holland, Amsterdam, 2003.

- [CDK<sup>+</sup>03b] Eric Cancès, Mireille Defranceschi, Werner Kutzelnigg, Claude Le Bris, and Yvon Maday. Computational quantum chemistry: a primer. In *Handbook of numerical analysis, Vol. X*, Handb. Numer. Anal., X, pages 3–270. North-Holland, Amsterdam, 2003.
- [CEL84] M. G. Crandall, L. C. Evans, and P.-L. Lions. Some properties of viscosity solutions of Hamilton-Jacobi equations. *Trans. Amer. Math. Soc.*, 282(2):487–502, 1984.
- [CLS07] Cancès, Eric , Legoll, Frédéric , and Stoltz, Gabriel . Theoretical and numerical comparison of some sampling methods for molecular dynamics. *Mathematical Modelling and Numerical Analysis*, 41(2):351–389, mar 2007.
- [CS01] B. Cano and A. M. Stuart. Underresolved Simulations of Heat Baths. *Journal of Computational Physics*, 169:193–214, May 2001.
- [CSS08] Carlsson, Jesper , Sandberg, Mattias , and Szepessy, Anders . Symplectic pontryagin approximations for optimal design. *ESAIM: Mathematical Modelling and Numerical Analysis*, PREPRINT, 2008.
- [DGL83] D. Dürr, S. Goldstein, and J. L. Lebowitz. A mechanical model for the Brownian motion of a convex body. *Probability Theory and Related Fields*, 62(4):427–448, 1983.
- [DGL81] D. Dürr, S. Goldstein, and J. L. Lebowitz. A mechanical model of Brownian motion. *Comm. Math. Phys.*, 78(4):507–530, 1980/81.
- [Dup94] Bruno Dupire. Pricing with a smile. *Risk*, 7(1):18–20, 1994.
- [Dur96] Richard Durrett. *Probability: theory and examples*. Duxbury Press, Belmont, CA, second edition, 1996.
- [EEHJ96] K. Eriksson, D. Estep, P. Hansbo, and C. Johnson. *Computational differential equations*. Cambridge University Press, Cambridge, 1996.
- [EHN96] Heinz W. Engl, Martin Hanke, and Andreas Neubauer. *Regularization of inverse problems*, volume 375 of *Mathematics and its Applications*. Kluwer Academic Publishers Group, Dordrecht, 1996.
- [Ein05] A. Einstein. Über die von der molekularkinetischen Theorie der Wärme geforderte Bewegung von in ruhenden Flüssigkeiten suspendierten Teilchen. *Annalen der Physik*, 322:549–560, 1905.
- [Evaa] L. C. Evans. An introduction to mathematical optimal control theory. <http://math.berkeley.edu/~evans/>.
- [Evab] L. C. Evans. An introduction to stochastic differential equations. <http://math.berkeley.edu/~evans/>.

- [Eva98] Lawrence C. Evans. *Partial differential equations*, volume 19 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 1998.
- [Fey98] Richard P. Feynman. *Statistical Mechanics: A Set of Lectures (Advanced Book Classics)*. Perseus Books Group, March 1998.
- [FHHØ08] S. Fournais, M. Hoffmann-Ostenhof, T. Hoffmann-Ostenhof, and T. Østergaard Sørensen. Analytic structure of many-body Coulombic wave functions. *ArXiv e-prints*, June 2008.
- [FK87] G. W. Ford and M. Kac. On the quantum langevin equation. *Journal of Statistical Physics*, 46:803–810, March 1987.
- [FKM65] G. W. Ford, M. Kac, and P. Mazur. Statistical Mechanics of Assemblies of Coupled Oscillators. *Journal of Mathematical Physics*, 6:504–515, April 1965.
- [Fre02] Berend Frenkel, Daanand Smit. *Understanding molecular simulation : from algorithms to applications*. Academic, San Diego, Calif., 2. ed. edition, 2002.
- [Gar91] C. W. Gardiner. *Quantum noise*, volume 56 of *Springer Series in Synergetics*. Springer-Verlag, Berlin, 1991.
- [Hag80] George A. Hagedorn. A time dependent Born-Oppenheimer approximation. *Commun. Math. Phys*, 77:1–19, 1980.
- [Hel88] Bernard Helffer. *Semi-classical analysis for the Schrödinger operator and applications*, volume 1336 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, 1988.
- [HK02] O. H. Hald and R. Kupferman. Asymptotic and numerical analyses for mechanical models of heat baths. *Journal of Statistical Physics*, 106(5):1121–1184, 2002.
- [HLW02] Ernst Hairer, Christian Lubich, and Gerhard Wanner. *Geometric numerical integration*, volume 31 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, 2002. Structure-preserving algorithms for ordinary differential equations.
- [Hul97] John C. Hull. *Options, futures, and other derivatives*. Prentice Hall International, London, 3. ed. edition, 1997.
- [Joh87] Claes Johnson. *Numerical solution of partial differential equations by the finite element method*. Studentlitteratur, Lund, 1987.
- [Kad00] Leo P. Kadanoff. *Statistical Physics: Statics, Dynamics and Renormalization*. World Scientific, 2000.

- [Koz00] V. V. Kozlov. Thermodynamics of Hamiltonian systems and the Gibbs distribution. *Dokl. Akad. Nauk*, 370(3):325–327, 2000.
- [KP92] Peter E. Kloeden and Eckhard Platen. *Numerical solution of stochastic differential equations*, volume 23 of *Applications of Mathematics (New York)*. Springer-Verlag, Berlin, 1992.
- [KS91] Ioannis Karatzas and Steven E. Shreve. *Brownian motion and stochastic calculus*, volume 113 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, second edition, 1991.
- [KS04] R. Kupferman and A. M. Stuart. Fitting SDE models to nonlinear Kac-Zwanzig heat bath models. *Phys. D*, 199(3-4):279–316, 2004.
- [KSTT02] R. Kupferman, A. M. Stuart, J. R. Terry, and P. F. Tupper. Long-term behaviour of large mechanical systems with random initial data. *Stoch. Dyn.*, 2(4):533–562, 2002.
- [LB05] Claude Le Bris. Computational chemistry from the perspective of numerical analysis. *Acta Numer.*, 14:363–444, 2005.
- [LG97] Don S. Lemons and Anthony Gythiel. Paul langevin’s 1908 paper “on the theory of brownian motion” [“sur la th[e-acute]orie du mouvement brownien,” c. r. acad. sci. (paris) [bold 146], 530–533 (1908)]. *American Journal of Physics*, 65(11):1079–1081, 1997.
- [LTSF05] X. Li, J. C. Tully, H. B. Schlegel, and M. J. Frisch. Ab initio ehrenfest dynamics. *Journal of Chemical Physics*, 123(8):1–7, 2005.
- [Mad26] E. Madelung. Quantum theory in hydrodynamic form. *Z. Phys.*, (40):322, 1926.
- [MH00] D. Marx and J. Hutter. Ab initio molecular dynamics: Theory and implementation. In J. Grotendorst, editor, *Modern Methods and Algorithms of Quantum Chemistry*, volume 1 of *NIC Series*, pages 301–449. John von Neumann Institute for Computing, Jülich, 2000.
- [Mil95] G. N. Milstein. *Numerical integration of stochastic differential equations*, volume 313 of *Mathematics and its Applications*. Kluwer Academic Publishers Group, Dordrecht, 1995. Translated and revised from the 1988 Russian original.
- [Mot08] N. F. Mott. On the theory of excitation by collision with heavy particles. *Mathematical Proceedings of the Cambridge Philosophical Society*, 27(04):553–560, 2008.
- [Øks98] Bernt Øksendal. *Stochastic differential equations*. Universitext. Springer-Verlag, Berlin, fifth edition, 1998. An introduction with applications.



- [PBG64] L. S. Pontryagin, V. G. Boltyanskii, R. V. Gamkrelidze, and E. F. Mishchenko. *The mathematical theory of optimal processes*. Translated by D. E. Brown. A Pergamon Press Book. The Macmillan Co., New York, 1964.
- [Ped99] Pablo Pedregal. Optimization, relaxation and Young measures. *Bull. Amer. Math. Soc. (N.S.)*, 36(1):27–58, 1999.
- [Pir84] Olivier Pironneau. *Optimal shape design for elliptic systems*. Springer Series in Computational Physics. Springer-Verlag, New York, 1984.
- [PST07] Gianluca Panati, Herbert Spohn, and Stefan Teufel. The time-dependent Born-Oppenheimer approximation. *M2AN Math. Model. Numer. Anal.*, 41(2):297–314, 2007.
- [San08] Mattias Sandberg. Convergence rates for an optimally controlled ginzburg-landau equation, 2008.
- [Sch28] Erwin Schrödinger. *Collected papers on Wave Mechanics*. Blackie and Son, London, 1928.
- [SS06] M. Sandberg and A. Szepessy. Convergence rates of symplectic Pontryagin approximations in optimal control theory. *M2AN*, 40(1), 2006.
- [Str86] Gilbert Strang. *Introduction to applied mathematics*. Wellesley-Cambridge Press, Wellesley, MA, 1986.
- [Sub95] Andreï I. Subbotin. *Generalized solutions of first-order PDEs*. Systems & Control: Foundations & Applications. Birkhäuser Boston Inc., Boston, MA, 1995. The dynamical optimization perspective, Translated from the Russian.
- [SW99] A. M. Stuart and J. O. Warren. Analysis and experiments for a computational model of a heat bath. *Journal of Statistical Physics*, 97(3):687–723, 1999.
- [Tan06] D. J. Tanner. *Introduction to Quantum Mechanics: A Time-dependent Perspective*. University Science Books, 2006.
- [Tul98a] J. C. Tully. Mixed quantum-classical dynamics. *Faraday Discuss.*, (110):407–419, 1998.
- [Tul98b] J. C. Tully. *Modern Methods for Multidimensional Dynamics Computation in Chemistry*, chapter 2. World Scientific Singapore, 1998.
- [Vog02] Curtis R. Vogel. *Computational methods for inverse problems*, volume 23 of *Frontiers in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002. With a foreword by H. T. Banks.

- [WD95] Sam Wilmott, Pauland Howison and Jeff Dewynne. *The mathematics of financial derivatives : a student introduction*. Cambridge Univ. Press, Cambridge, 1995.
- [You69] L. C. Young. *Lectures on the calculus of variations and optimal control theory*. Foreword by Wendell H. Fleming. W. B. Saunders Co., Philadelphia, 1969.
- [Zwa73] R. Zwanzig. Nonlinear generalized Langevin equations. *J. Stat. Phys.*, 9:215–220, 1973.