

# Lecture notes for Math 227

Lenya Ryzhik\*

February 19, 2013

These notes come from the original notes by Jim Nolen, as well as other sources.

## 1 The Haar functions and the Brownian motion

### Prologue: a probabilistic interpretation for a difference equation

A good way to understand how the probabilistic interpretation of some PDEs comes about, is to start with the discrete equations. Consider the finite difference analog of the Laplace equation:

$$u(x+1, y) + u(x-1, y) + u(x, y+1) + u(x, y-1) - 4u(x, y) = 0, \quad (1.1)$$

which is the discrete analog of the Laplace equation

$$-\Delta u = 0,$$

where

$$\Delta u = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2},$$

in two dimensions. In order to formulate a (discrete) boundary value problem, let  $U$  be a domain of the two-dimensional square lattice  $\mathbb{Z}^2$ , and let  $u(x, y)$  solve the difference equation (1.1), with the boundary condition  $u(x, y) = g(x, y)$  on the boundary  $\partial U$ . Here  $g(x, y)$  is a prescribed non-negative function, which is positive somewhere.

We claim that the solution of this problem has the following probabilistic interpretation. Let  $(X(t), Y(t))$  be the standard random walk on the lattice  $\mathbb{Z}^2$  – the probability to go up down, left or right is equal to 1/4, and let it start at the point  $(x, y)$ :  $X(0) = x$ ,  $Y(0) = y$ . Let  $(\bar{x}, \bar{y})$  be the first point where  $(X(t), Y(t))$  reaches the boundary  $\partial U$  of the domain. The point  $(\bar{x}, \bar{y})$  is, of course, random. The beautiful observation is that the function  $v(x, y) = \mathbb{E}(g(\bar{x}, \bar{y}))$  gives a solution of (1.1), connecting this discrete problem to the random walk. Why? First, it is immediate that if  $(x, y)$  is on the boundary then, of course,  $\bar{x} = x$  and  $\bar{y} = y$ , so  $v(x, y) = g(x, y)$  in that case. On the other hand, if  $(x, y)$  is inside  $U$  then

$$v(x, y) = \frac{1}{4}(v(x+1, y) + v(x-1, y) + v(x, y+1) + v(x, y-1))$$

---

\*Department of Mathematics, Stanford University, Stanford CA 94305; ryzhik@math.stanford.edu

simply from the definition of the random walk, and the definition of  $v(x, y)$  – we can go in any of the four possible directions with the probability equal to  $1/4$  and then we start afresh.

Now, if we let the mesh size be not 1 but  $h > 0$  and let  $h \downarrow 0$ , the discrete equation (1.1) becomes the Laplace equation, while the random walk becomes the Brownian motion. More precisely, solution of the boundary value problem

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0, \quad (x, y) \in U, \quad (1.2)$$

with the boundary condition  $u(x, y) = g(x, y)$  for  $(x, y) \in \partial U$ , has the following probabilistic interpretation. Consider a Brownian motion  $B(t; x, y)$  that starts at a point  $(x, y) \in U$  and let  $(\bar{x}, \bar{y})$  be a (random) point where  $B(t; x, y)$  hits the boundary  $\partial U$  for the first time. Then solution of (1.2) is  $u(x, y) = \mathbb{E}(g(\bar{x}, \bar{y}))$ .

From the heuristic point of view, we can deduce immediately some properties of the Laplace equation. For example, if  $g(x, y) > 0$  in some open subset  $S \subset U$ , and  $g(x, y) \geq 0$  for all  $(x, y) \in \partial U$ , then with a positive probability we have  $(\bar{x}, \bar{y}) \in S$  for all  $(x, y) \in U$  (this is not really obvious but is true), which means that  $u(x, y) = \mathbb{E}(g(\bar{x}, \bar{y})) > 0$  for all  $(x, y) \in U$ .

It is also easy to deduce the maximum principle from the probabilistic interpretation: it is easy to see that  $\mathbb{E}(g(\bar{x}, \bar{y})) \leq \sup_{(x, y) \in \partial U} g(x, y)$  – expected value of a function can not exceed its maximum.

In order to really use such ideas we need some probabilistic background such as what is the Brownian motion, and this is how we will start – with a construction of the Brownian motion on the real line. This will be done using the Haar functions in a very explicit way. We will rarely prove purely probabilistic results in this course but this construction is so basic and explicit that it would be criminal to omit it.

## 1.1 The Haar functions and their completeness

### The Haar functions

The basic Haar function is

$$\psi(x) = \begin{cases} 1 & \text{if } 0 \leq x < 1/2, \\ -1 & \text{if } 1/2 \leq x < 1, \\ 0 & \text{otherwise.} \end{cases} \quad (1.3)$$

It has mean zero

$$\int_0^1 \psi(x) dx = 0,$$

and is normalized so that

$$\int_0^1 \psi^2(x) dx = 1.$$

The rescaled and shifted Haar functions are

$$\psi_{jk}(x) = 2^{j/2} \psi(2^j x - k), \quad j, k \in \mathbb{Z}.$$

These functions form an orthonormal set in  $L^2(\mathbb{R})$  because if  $j = j'$  and  $k \neq k'$  then

$$\int_{\mathbb{R}} \psi_{jk}(x)\psi_{j'k'}(x)dx = 2^j \int_{\mathbb{R}} \psi(2^j x - k)\psi(2^j x - k')dx = 0$$

because  $\psi(y - k)\psi(y - k') = 0$  for any  $y \in \mathbb{R}$  and  $k \neq k'$ . On the other hand, if  $j \neq j'$ , say,  $j < j'$ , then

$$\begin{aligned} \int_{\mathbb{R}} \psi_{jk}(x)\psi_{j'k'}(x)dx &= 2^{j/2+j'/2} \int_{\mathbb{R}} \psi(2^j x - k)\psi(2^{j'} x - k')dx \\ &= 2^{j'/2-j/2} \int_{\mathbb{R}} \psi(y)\psi(2^{j'-j}y + 2^{j'-j}k - k')dy \\ &= 2^{j'/2-j/2} \int_0^{1/2} \psi(2^{j'-j}y + 2^{j'-j}k - k')dy - 2^{j'/2-j/2} \int_{1/2}^1 \psi(2^{j'-j}y + 2^{j'-j}k - k')dy. \end{aligned}$$

Both of the integrals above equal to zero. Indeed,  $2^{j'-j} \geq 2$ , hence, for instance,

$$\int_0^{1/2} \psi(2^{j'-j}y + 2^{j'-j}k - k')dy = 2^{j-j'} \int_0^{2^{j'-j-1}} \psi(y + 2^{j'-j}k - k')dy = 0,$$

because

$$\int_m^n \psi(y)dy = 0,$$

for all  $m, n \in \mathbb{Z}$ , and  $j' > j$ . Finally, when  $j = j'$ ,  $k = k'$  we have

$$\int_{\mathbb{R}} |\psi_{jk}(x)|^2 = 2^j \int_{\mathbb{R}} |\psi(2^j x - k)|^2 dx = \int_{\mathbb{R}} |\psi(x - k)|^2 dx = 1.$$

The Haar coefficients of a function  $f \in L^2(\mathbb{R})$  are defined as the inner products

$$c_{jk} = \int f(x)\psi_{jk}(x)dx, \tag{1.4}$$

and the Haar series of  $f$  is

$$\sum_{j,k \in \mathbb{Z}} c_{jk}\psi_{jk}(x). \tag{1.5}$$

Orthonormality of the family  $\{\psi_{jk}\}$  ensures that

$$\sum_{j,k} |c_{jk}|^2 \leq \|f\|_{L^2}^2 < +\infty,$$

and the series (1.5) converges in  $L^2(\mathbb{R})$ . In order to show that it actually converges to the function  $f$  itself we need to prove that the Haar functions form a basis for  $L^2(\mathbb{R})$ .

## Completeness of the Haar functions

To show that Haar functions form a basis in  $L^2(\mathbb{R})$  we consider the dyadic projections  $P_n$  defined as follows. Given  $f \in L^2(\mathbb{R})$ , and  $n, k \in \mathbb{Z}$ , consider the intervals

$$I_{nk} = ((k-1)/2^n, k/2^n],$$

then

$$P_n f(x) = \int_{I_{nk}} f dx = 2^n \int_{I_{nk}} f dx, \quad \text{for } x \in I_{nk}.$$

The function  $P_n f$  is constant on each of the dyadic intervals  $I_{nk}$ . In particular, each Haar function  $\psi_{jk}$  satisfies  $P_n \psi_{jk}(x) = 0$  for  $j \geq n$ , while  $P_n \psi_{jk}(x) = \psi_{jk}(x)$  for  $j < n$ . We claim that, actually, for any  $f \in L^2(\mathbb{R})$  we have

$$P_{n+1}f - P_n f = \sum_{k \in \mathbb{Z}} c_{nk} \psi_{nk}(x), \quad (1.6)$$

with the Haar coefficients  $c_{nk}$  given by (1.4). Indeed, let  $x \in I_{nk}$  and write

$$I_{nk} = \left( \frac{(k-1)}{2^n}, \frac{k}{2^n} \right] = \left( \frac{2(k-1)}{2^{n+1}}, \frac{2k-1}{2^{n+1}} \right] \cup \left( \frac{2k-1}{2^{n+1}}, \frac{2k}{2^{n+1}} \right] = I_{n+1,2k-1} \cup I_{n+1,2k}.$$

The function  $P_n f$  is constant on the whole interval  $I_{nk}$  while  $P_{n+1}f$  is constant on each of the sub-intervals  $I_{n+1,2k-1}$  and  $I_{n+1,2k}$ . In addition,

$$\int_{I_{nk}} (P_n f) dx = \int_{I_{nk}} (P_{n+1}f) dx.$$

This means exactly that

$$P_{n+1}(x) = P_n f(x) + c_{nk} \psi_{nk}(x) \text{ for } x \in I_{nk},$$

which is (1.6).

As a consequence of (1.6) we deduce that

$$P_{n+1}f(x) - P_m f(x) = \sum_{j=-m}^n \sum_{k \in \mathbb{Z}} c_{jk} \psi_{jk}(x), \quad (1.7)$$

for all  $m, n \in \mathbb{Z}$  with  $n > m$ . It remains to show that for any  $f \in L^2(\mathbb{R})$  we have

$$\lim_{m \rightarrow +\infty} P_m f(x) = 0, \quad \lim_{n \rightarrow +\infty} P_n f(x) = f(x), \quad (1.8)$$

both in the  $L^2$ -sense. The operators  $P_n f$  are uniformly bounded because for all  $n, k \in \mathbb{Z}$  we have

$$\int_{I_{nk}} |(P_n f)(x)|^2 dx = 2^{-n} 2^{2n} \left| \int_{I_{nk}} f(y) dy \right|^2 \leq \int_{I_{nk}} |f(y)|^2 dy.$$

Summing over  $k \in \mathbb{Z}$  for a fixed  $n$  we get

$$\int_{\mathbb{R}} |P_n f(x)|^2 \leq \int_{\mathbb{R}} |f(x)|^2,$$

thus  $\|P_n f\|_{L^2} \leq \|f\|_{L^2}$ . Uniform boundedness of  $P_n$  implies that it is sufficient to establish both limits in (1.8) for functions  $f \in C_c(\mathbb{R})$ . However, for such  $f$  we have, on one hand,

$$|P_{-m} f(x)| \leq \frac{1}{2^m} \int_{\mathbb{R}} |f(x)| dx \rightarrow 0 \text{ as } m \rightarrow +\infty,$$

and, on the other,  $f$  is uniformly continuous on  $\mathbb{R}$ , so that  $\|P_n f(x) - f(x)\|_{L^\infty} \rightarrow 0$  as  $n \rightarrow +\infty$ , which, as both  $P_n f$  and  $f$  are compactly supported, implies the second limit in (1.8). Therefore,  $\psi_{jk}$  form an orthonormal basis in  $L^2(\mathbb{R})$  and every function  $f \in L^2(\mathbb{R})$  has the representation

$$f(x) = \sum_{j,k=-\infty}^{\infty} c_{jk} \psi_{jk}(x), \quad c_{jk} = \int_{\mathbb{R}} f(y) \psi_{jk}(y) dy. \quad (1.9)$$

### Possible lack of convergence in $L^1(\mathbb{R})$

Let us note that the Haar series need not converge in  $L^1(\mathbb{R})$ . In particular, it can not be integrated term-wise to conclude that

$$\int_{\mathbb{R}} f(x) dx = 0,$$

which obviously can not be true for all  $f \in L^2(\mathbb{R})$ . Consider, for example, the function  $f(x) = \chi_{[0,1]}(x)$ , that is,  $f(x) = 1$  for  $x \in [0, 1]$  and  $f(x) = 0$  otherwise. Its Haar coefficients are

$$c_{jk} = \int_0^1 \psi_{jk}(x) dx = 2^{j/2} \int_0^1 \psi(2^j x - k) dx = 2^{-j/2} \int_0^{2^j} \psi(x - k) dx = 2^{-j/2} \int_{-k}^{-k+2^j} \psi(x) dx.$$

We see immediately that  $c_{jk} = 0$  for  $k \neq 0$ , and also for  $j \geq 0$ . On the other hand, for  $j < 0$  and  $k = 0$  we have  $c_{j,0} = 2^{j/2}$ . Therefore, the Haar series for the function  $f$  is

$$\sum_{j < 0} 2^j \psi(2^j x) = \sum_{j > 0} \frac{1}{2^j} \psi\left(\frac{x}{2^j}\right).$$

The partial sums

$$S_N(x) = \sum_{j=1}^N \frac{1}{2^j} \psi\left(\frac{x}{2^j}\right)$$

look as follows:

$$S_N(x) = \sum_{j=1}^N \frac{1}{2^j} = 1 - \frac{1}{2^N}, \quad \text{for } x \in [0, 1],$$

and,

$$S_N(x) = \sum_{j=2}^N \frac{1}{2^j} - \frac{1}{2} = -\frac{1}{2^N} \quad \text{for } x \in [1, 2^N],$$

while  $S_N(x) = 0$  for  $x > 2^N$ . The function  $S_N(x)$  has tails that are not uniformly integrable, and the series does not converge to  $f(x) = \chi_{[0,1]}(x)$  in  $L^1$ . However, we have

$$\int_{\mathbb{R}} |S_N(x) - 1|^2 dx = \frac{1}{2^N} \rightarrow 0,$$

hence the series converges in  $L^2(\mathbb{R})$ , as it should. This is a general phenomenon – the Haar series converges in  $L^p$  for  $1 < p < \infty$  but not in  $L^1$  since it has fat tails that decay only as  $1/x$ , which is not sufficient for the  $L^1$ -convergence.

## 1.2 The Brownian motion

Brownian motion is a random process  $X_t(\omega)$ ,  $t \geq 0$  defined on a probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  which has the following properties:

- (i) The function  $X_t(\omega)$  is continuous in  $t$  for a.e. realization  $\omega$ .
- (ii) For all  $0 \leq s < t < +\infty$  the random variable  $X_t(\omega) - X_s(\omega)$  is Gaussian with mean zero and variance  $t - s$ :

$$\mathbb{E}(X(t) - X(s)) = 0, \quad \mathbb{E}(X(t) - X(s))^2 = t - s.$$

- (iii) For any subdivision  $0 = t_0 < t_1 < \dots < t_N = t$  of the interval  $[0, t]$ , the random variables  $X_{t_1} - X_{t_0}, \dots, X_{t_N} - X_{t_{N-1}}$  are independent.

### Construction of the Brownian motion

We will construct the Brownian motion on the interval  $0 \leq t \leq 1$  – the restriction to a finite interval is a simple convenience but by no means a necessity. The Haar functions  $\psi_{jk}(x)$ , with  $j \geq 0$ ,  $0 \leq k \leq 2^j - 1$ , form a basis for the space  $L^2[0, 1]$ . Let us denote accordingly  $\phi_n(x) = \psi_{jk}(x)$  for  $n = 2^j + k$ ,  $0 \leq k \leq 2^j - 1$ , and  $\phi_0(x) = 1$  so that  $\{\phi_n\}$  form an orthonormal basis for  $L^2[0, 1]$ . Let  $Z_n(\omega)$ ,  $n \geq 0$ , be a collection of independent Gaussian random variables of mean zero and variance one, that is,

$$P(Z_n < y) = \int_{-\infty}^y e^{-y^2/2} \frac{dy}{\sqrt{2\pi}}.$$

We will show that the process

$$X_t(\omega) = \sum_{n=0}^{\infty} Z_n(\omega) \int_0^t \phi_n(s) ds \tag{1.10}$$

is a Brownian motion.

First, we need to verify that the series (1.10) converges in  $L^2(\Omega)$  for a fixed  $t \in [0, 1]$ . Note that

$$\mathbb{E} \left( \sum_{k=n}^m Z_k(\omega) \int_0^t \phi_k(s) ds \right)^2 = \sum_{k=n}^m \left( \int_0^t \phi_k(s) ds \right)^2 = \sum_{k=n}^m \langle \chi_{[0,t]}, \phi_k \rangle^2.$$

As  $\phi_k$  form a basis for  $L^2[0, 1]$ , it follows that the series (1.10) satisfies the Cauchy criterion and thus converges in  $L^2(\Omega)$ . Moreover, for any  $0 \leq s < t \leq 1$  we have

$$\begin{aligned} \mathbb{E} (X_t - X_s)^2 &= \mathbb{E} \left( \sum_{k=0}^{\infty} Z_k(\omega) \int_s^t \phi_k(u) du \right)^2 = \sum_{k=0}^{\infty} \left( \int_s^t \phi_k(u) du \right)^2 = \sum_{k=0}^{\infty} \langle \chi_{[s,t]}, \phi_k \rangle^2 \\ &= \|\chi_{[s,t]}\|_{L^2}^2 = t - s, \end{aligned}$$

hence the increments  $X_t - X_s$  have the correct variance. Let us show that they are independent: for  $0 \leq t_0 < t_1 \leq t_2 < t_3 \leq 1$ :

$$\begin{aligned} \mathbb{E} ((X_{t_3} - X_{t_2})(X_{t_1} - X_{t_0})) &= \mathbb{E} \left( \sum_{k=0}^{\infty} \int_{t_2}^{t_3} \phi_k(u) du \int_{t_0}^{t_1} \phi_k(u') du' \right) \\ &= \sum_{k=0}^{\infty} \langle \chi_{[t_2,t_3]}, \phi_k \rangle \langle \chi_{[t_0,t_1]}, \phi_k \rangle = \langle \chi_{[t_2,t_3]}, \chi_{[t_0,t_1]} \rangle = 0. \end{aligned}$$

As the variables  $X_t - X_s$  are jointly Gaussian, independence of the increments follows.

### Continuity of the Brownian motion

In order to prove continuity of the process  $X_t(\omega)$  defined by the series (1.10) we show that the series converges uniformly in  $t$  almost surely in  $\omega$ . To this end let us show that

$$M(\omega) = \sup_n \frac{|Z_n(\omega)|}{\sqrt{\log n}} < +\infty \text{ almost surely in } \omega. \quad (1.11)$$

Note that, for each  $n \geq 0$ :

$$\mathbb{P} \left( |Z_n(\omega)| \geq 2\sqrt{\log n} \right) \leq e^{-(2\sqrt{\log n})^2/2} = \frac{1}{n^2},$$

thus

$$\sum_{n=0}^{\infty} \mathbb{P} \left( |Z_n(\omega)| \geq 2\sqrt{\log n} \right) < +\infty.$$

The Borel-Cantelli lemma implies that almost surely the event  $\{|Z_n(\omega)| \geq 2\sqrt{\log n}\}$  happens only finitely many times, so that  $|Z_n(\omega)| < 2\sqrt{\log n}$  for all  $n \geq n_0(\omega)$  almost surely, and (1.11) follows.

Another useful observation is that for each fixed  $t \geq 0$  and  $j \in \mathbb{N}$  there exists only one  $k$  so that

$$\int_0^t \phi_{2^j+k}(s) ds \neq 0,$$

and for that  $k$  we have

$$\left| \int_0^t \phi_{2^j+k}(s) ds \right| \leq 2^{j/2} 2^{-j} = \frac{1}{2^{j/2}}.$$

Hence, we may estimate the dyadic blocs, using (1.11):

$$\left| \sum_{k=0}^{2^j-1} Z_{2^j+k}(\omega) \int_0^t \phi_{2^j+k}(s) ds \right| \leq M(\omega) \sqrt{(j+1) \log 2} \sum_{k=0}^{2^j-1} \left| \int_0^t \psi_{jk}(s) ds \right| \leq \frac{\sqrt{j} M_1(\omega)}{2^{j/2}}.$$

Therefore, the dyadic blocs are bounded by a convergent series which does not depend on  $t \in [0, 1]$ , hence the sum  $X_t(\omega)$  of the series is a continuous function for a.e.  $\omega$ .

### Nowhere differentiability of the Brownian motion

**Theorem 1.1** *The Brownian path  $X_t(\omega)$  is nowhere differentiable for almost every  $\omega$ .*

**Proof.** Let us fix  $\beta > 0$ . Then if  $\dot{X}_s$  exists at some  $s \in [0, 1]$  and  $|\dot{X}_s| < \beta$  then there exists  $n_0$  so that

$$|X_t - X_s| \leq 2\beta|t - s| \text{ if } |t - s| \leq \frac{2}{n} \quad (1.12)$$

for all  $n > n_0$ . Let  $A_n$  be the set of functions  $x(t) \in C[0, 1]$  for which (1.12) holds for some  $s \in [0, 1]$ . Then  $A_n \subset A_{n+1}$  and the set  $A = \bigcup_{n=1}^{\infty} A_n$  includes all functions  $x(t) \in C[0, 1]$  such that  $|\dot{x}(s)| \leq \beta$  at some point  $s \in [0, 1]$ .

The next step is to replace (1.12) by a discrete set of conditions – this is a standard trick in such situations. Assume that (1.12) holds for a function  $x(t) \in C[0, 1]$  and let  $k = \sup\{j : j/n \leq s\}$ , then

$$y_k = \max\left(\left|x\left(\frac{k+2}{n}\right) - x\left(\frac{k+1}{n}\right)\right|, \left|x\left(\frac{k+1}{n}\right) - x\left(\frac{k}{n}\right)\right|, \left|x\left(\frac{k}{n}\right) - x\left(\frac{k-1}{n}\right)\right|\right) \leq \frac{8\beta}{n}.$$

Therefore, if we denote by  $B_n$  the set of all functions  $x(t) \in C[0, 1]$  for which  $y_k \leq 8\beta/n$  for some  $k$ , then  $A_n \subseteq B_n$ . Therefore, in order to show that  $\mathbb{P}(A) = 0$  it suffices to check that

$$\lim_{n \rightarrow \infty} \mathbb{P}(B_n) = 0. \quad (1.13)$$

This, however, can be estimated directly, using translation invariance of the Brownian motion:

$$\begin{aligned} \mathbb{P}(B_n) &\leq \sum_{k=1}^{n-2} \mathbb{P}\left[\max\left[\left|X\left(\frac{k+2}{n}\right) - X\left(\frac{k+1}{n}\right)\right|, \left|X\left(\frac{k+1}{n}\right) - X\left(\frac{k}{n}\right)\right|, \left|X\left(\frac{k}{n}\right) - X\left(\frac{k-1}{n}\right)\right|\right] \leq \frac{8\beta}{n}\right] \\ &\leq n\mathbb{P}\left[\max\left[\left|X\left(\frac{3}{n}\right) - X\left(\frac{2}{n}\right)\right|, \left|X\left(\frac{2}{n}\right) - X\left(\frac{1}{n}\right)\right|, \left|X\left(\frac{1}{n}\right)\right|\right] \leq \frac{8\beta}{n}\right] \\ &= n\mathbb{P}\left[\left|X\left(\frac{1}{n}\right)\right| \leq \frac{8\beta}{n}\right]^3 = n\left(\sqrt{\frac{n}{2\pi}} \int_{-8\beta/n}^{8\beta/n} e^{-nx^2/2} dx\right)^3 \leq n\left(\sqrt{\frac{n}{2\pi}} \frac{16\beta}{n}\right)^3 \leq \frac{C}{\sqrt{n}}, \end{aligned}$$

which implies (1.13). It follows that  $\mathbb{P}(A) = 0$  as well, hence Brownian motion is nowhere differentiable with probability one.  $\square$

**Corollary 1.2** *Brownian motion does not have bounded variation with probability one.*

## 2 Stochastic integration

We would like now to understand how to interpret and ODE of the form

$$\frac{dX}{dt} = \dot{B}(t).$$

The right side certainly does not make sense since  $B(t)$  is almost surely not differentiable in time but that should not stop us from trying. More generally, we would like to look at ODE's of the form

$$\frac{dX}{dt} = b(t, X) + \sigma(t, X)\dot{B}(t).$$

A natural approach is to start with a discretized version

$$X_{k+1} - X_k = b(t_k^*, X_k)\Delta t + \sigma(t_k^*, X_k)\Delta B_k, \quad (2.1)$$

where  $\Delta$  is the time step, and  $\Delta B_k = B(t_{k+1}) - B(t_k)$  is the increment of the Brownian motion. The time  $t_k^*$  is taken on the interval  $[t_k, t_{k+1}]$ . We take for the moment  $t_k^* = t_k$  - we will later discuss what happens with other choices of  $t_k^*$ . In the "integrated form" (2.1) becomes

$$X_k = X_0 + \sum_{j=0}^{k-1} b(t_j, X_j)\Delta t_j + \sum_{j=0}^{k-1} \sigma(t_j, X_j)\Delta B_j. \quad (2.2)$$

The question we need to understand is whether there exists a limit as  $\Delta t \rightarrow 0$  for the solutions of such discrete equations.

## The Ito integral

With the above goal in mind, we we would like to define

$$\int_S^T f(t, \omega) dB_t(\omega),$$

with  $0 < S < T$ . As in the definition of any integral, we start with simple functions of the form

$$\phi(t, \omega) = \sum_{j \geq 0} e_j(\omega) \chi_{[j/2^n, (j+1)/2^n)}(t).$$

The only reasonable definition of the integral for simple functions is to set

$$\int_S^T \phi(t, \omega) dB_t(\omega) = \sum_{j \geq 0} e_j(\omega) [B(t_{j+1}) - B(t_j)]. \quad (2.3)$$

In order to understand what can happen with this definition, let us consider the following two examples:

$$\phi_1(t, \omega) = \sum_{j \geq 0} B\left(\frac{j}{2^n}, \omega\right) \chi_{[j/2^n, (j+1)/2^n)}(t),$$

and

$$\phi_2(t, \omega) = \sum_{j \geq 0} B\left(\frac{j+1}{2^n}, \omega\right) \chi_{[j/2^n, (j+1)/2^n)}(t).$$

Both of these functions are, supposedly, approximating  $\phi(t, \omega) = B_t(\omega)$ , so their integrals should have the same limit as  $n \rightarrow +\infty$ . Let us see (we assume for simplicity that  $T = M2^{-n}$  is an integer multiple of  $2^{-n}$ ):

$$\int_0^T \phi_1(t, \omega) dB_t(\omega) = \sum_{j=0}^{M-1} B\left(\frac{j}{2^n}\right) [B\left(\frac{j+1}{2^n}\right) - B\left(\frac{j}{2^n}\right)],$$

and

$$\int_0^T \phi_2(t, \omega) dB_t(\omega) = \sum_{j=0}^M B\left(\frac{j+1}{2^n}\right) \left[ B\left(\frac{j+1}{2^n}\right) - B\left(\frac{j}{2^n}\right) \right].$$

Hence,  $\phi_1(t, \omega)$  corresponds to approximating  $\phi(t, \omega)$  by taking  $t_j^* = t_j$  while  $\phi_2(t, \omega)$  corresponds to  $t_j^* = t_{j+1}$ . We compute:

$$\mathbb{E} \left( \int_0^T \phi_1(t, \omega) dB_t(\omega) \right) = \sum_{j=0}^M \mathbb{E} \left( B\left(\frac{j}{2^n}\right) \left[ B\left(\frac{j+1}{2^n}\right) - B\left(\frac{j}{2^n}\right) \right] \right) = 0,$$

while

$$\begin{aligned} \mathbb{E} \left( \int_0^T \phi_2(t, \omega) dB_t(\omega) \right) &= \sum_{j=0}^M \mathbb{E} \left( B\left(\frac{j+1}{2^n}\right) \left[ B\left(\frac{j+1}{2^n}\right) - B\left(\frac{j}{2^n}\right) \right] \right) \\ &= \sum_{j=0}^M \mathbb{E} \left( \left( B\left(\frac{j+1}{2^n}\right) - B\left(\frac{j}{2^n}\right) \right) \left[ B\left(\frac{j+1}{2^n}\right) - B\left(\frac{j}{2^n}\right) \right] \right) = \sum_{j=0}^M \frac{1}{2^n} = T. \end{aligned}$$

Therefore, the "integrals" of  $\phi_1(t, \omega)$  and  $\phi_2(t, \omega)$  differ in a non-trivial way that does not vanish as  $n \rightarrow +\infty$  – the choice of  $t_j^*$  matters! There are two canonical choices:  $t_j^* = t_j$  gives rise to the Ito integral, while  $t_j^* = (t_j + t_{j+1})/2$  leads to the Stratonovich integral.

## Integrable functions

We begin with the following definition.

**Definition 2.1** *Let  $\{\mathcal{N}_t\}$  be an increasing family of  $\sigma$ -algebras. A process  $g(t, \omega)$  is  $\mathcal{N}_t$ -adapted if  $g(t, \omega)$  is  $\mathcal{N}_t$ -measurable for each  $t \geq 0$ .*

A typical situation when this definition is used is when  $\mathcal{N}_t$  is generated by a process  $X(t, \omega)$  –  $\mathcal{N}_t$  is the collection of all events that depend on  $X(s, \omega)$  for  $0 \leq s \leq t$  but not on  $X(s, \omega)$  for  $s > t$ . In turn,  $g(t, \omega)$  is  $\mathcal{N}_t$ -adapted if  $g(t, \omega)$  depends only on  $X(s, \omega)$  for  $0 \leq s \leq t$ . For example,

$$g(t, \omega) = \int_0^t X(s, \omega) ds$$

is  $\mathcal{N}_t$ -adapted, while  $g(t, \omega) = \max_{t \leq s \leq t+1} |X(s, \omega)|$  is not  $\mathcal{N}_t$ -adapted.

We will define the Ito integral for functions  $f(t, \omega)$  (measurable in both variables) for which the following two conditions hold, on a time interval  $[0, T]$ :

(i)  $f(t, \omega)$  is  $\mathcal{F}_t$ -adapted (here  $\mathcal{F}_t$  is the  $\sigma$ -algebra of events generated by  $\{B_s : 0 \leq s \leq t\}$ , and

$$(ii) \quad \mathbb{E} \left( \int_0^T f(t, \omega)^2 dt \right) < +\infty.$$

We will denote by  $V$  the class of functions for which both conditions (i) and (ii) hold.

## Ito integral for elementary functions

**Definition 2.2** A function  $\phi \in V$  is elementary if it has the form

$$\phi(t, \omega) = \sum_j e_j(\omega) \chi_{[t_j, t_{j+1})}(t).$$

For an elementary function  $\phi(t, \omega)$  we set

$$\int_S^T \phi(t, \omega) dB_t(\omega) = \sum_{j \geq 0} e_j(\omega) (B(t_{j+1}) - B(t_j)).$$

If  $\phi \in V$  then it is piece-wise constant in time, and  $e_j(\omega)$  has to be  $\mathcal{F}_{t_j}$ -measurable – this follows from the fact that it is  $\mathcal{F}_t$ -adapted.

A very important observation is the following:

**Lemma 2.3** (Ito isometry) If  $\phi(t, \omega)$  is bounded and elementary then

$$\mathbb{E} \left( \int_S^T \phi(t, \omega) dB_t(\omega) \right)^2 = \mathbb{E} \left[ \int_S^T \phi^2(t, \omega) dt \right]. \quad (2.4)$$

**Proof.** Let us take

$$\phi(t, \omega) = \sum_j e_j(\omega) \chi_{[t_j, t_{j+1})}(t),$$

and compute ( $\Delta B_j = B(t_{j+1}) - B(t_j)$ , and similarly for  $\Delta B_i$ ):

$$\mathbb{E} \left( \int_S^T \phi(t, \omega) dB_t(\omega) \right)^2 = \mathbb{E} \left( \sum_{i,j} e_i(\omega) e_j(\omega) \Delta B_i \Delta B_j \right). \quad (2.5)$$

Note that when  $i \neq j$ , say,  $i > j$ , then  $e_i(\omega)$ ,  $e_j(\omega)$  and  $\Delta B_j$  all depend only on the Brownian motion until the time  $t_i$ , and thus these quantities are independent from the forward increment  $B(t_{i+1}) - B(t_i)$ , whence

$$\mathbb{E}(e_i(\omega) e_j(\omega) \Delta B_i \Delta B_j) = \mathbb{E}(e_i(\omega) e_j(\omega) \Delta B_j) \mathbb{E}(\Delta B_i) = 0, \quad \text{for } i > j.$$

Using this in (2.5) gives

$$\mathbb{E} \left( \int_S^T \phi(t, \omega) dB_t(\omega) \right)^2 = \mathbb{E} \left( \sum_{i=j} e_i(\omega) e_j(\omega) \Delta B_i \Delta B_j \right) = \sum_j \mathbb{E}(e_j^2(\omega) (\Delta B_j)^2). \quad (2.6)$$

Once again,  $e_j(\omega)$  depends only on the Brownian motion until the time  $t_j$  and is, therefore, independent of the forward increment  $\Delta B(t_j)$ , giving

$$\mathbb{E}(e_j^2(\omega) (\Delta B_j)^2) = \mathbb{E}(e_j^2(\omega)) \mathbb{E}((\Delta B_j)^2) = \sum_j \mathbb{E}(e_j^2(\omega)) \Delta t_j.$$

We conclude that

$$\mathbb{E} \left( \int_S^T \phi(t, \omega) dB_t(\omega) \right)^2 = \sum_j \mathbb{E}(e_j^2(\omega)) \Delta t_j, \quad (2.7)$$

which proves Lemma 2.3.  $\square$

## Extension to non-elementary functions

We now extend, gingerly and slowly, the notion of the Ito integral to non-elementary functions. The idea is to show that bounded elementary functions are dense in  $V$  and then use Ito's isometry property. We will prove density of bounded elementary functions in  $V$  in three steps: (1) Show that bounded continuous functions can be approximated by bounded elementary functions, (2) Show that bounded functions in  $V$  can be approximated by bounded continuous functions, and (3) Show that bounded functions are dense in  $V$ .

**Step 1.** Let  $g \in V$  be bounded, and also continuous in  $t$  for each  $\omega$ . We claim that there exists a sequence of elementary functions  $\phi_n$  so that

$$\mathbb{E} \left( \int_S^T (g - \phi_n)^2 dt \right) \rightarrow 0 \text{ as } n \rightarrow +\infty. \quad (2.8)$$

This is true for the natural piece-wise constant approximation

$$\phi_n(t, \omega) = \sum_j g(t_j, \omega) \chi_{[t_j, t_{j+1})}(t).$$

To see that, note first that

$$\int_S^T (g - \phi_n)^2 ds \rightarrow 0 \text{ as } n \rightarrow +\infty, \quad (2.9)$$

for each realization  $\omega$  fixed, because  $g(t, \omega)$  is continuous for each  $\omega$  fixed. As the function  $g(t, \omega)$  is bounded, (2.9) and the Lebesgue dominated convergence theorem imply that

$$\mathbb{E} \int_S^T (g - \phi_n)^2 ds \rightarrow 0 \text{ as } n \rightarrow +\infty,$$

which is (2.8).

**Step 2.** Let  $h \in V$  be bounded. We claim that there exists a sequence of bounded continuous functions  $\psi_n$  so that

$$\mathbb{E} \left( \int_S^T (h - \psi_n)^2 dt \right) \rightarrow 0 \text{ as } n \rightarrow +\infty. \quad (2.10)$$

Indeed, suppose that  $|h(t, \omega)| \leq M$  and take a smooth function  $g_n(t) \geq 0$  such that  $g_n(t) = 0$  for  $t$  outside  $[-1/n, 0]$  and

$$\int_{\mathbb{R}} g_n(t) dt = 1.$$

Now, set

$$\psi_n(t) = \int_0^t g_n(s-t) h(s, \omega) ds,$$

then  $\psi_n(t, \omega)$  is continuous in  $t$ , and  $|\psi_n(t)| \leq M$ . Moreover, as  $g_n(t) = 0$  for  $t \geq 0$ , and  $h(t, \omega)$  is  $\mathcal{F}_t$ -adapted, the functions  $\psi_n(t, \omega)$  are also  $\mathcal{F}_t$ -adapted. It is also easy to check that

$$\int_S^T (h(t, \omega) - \psi_n(t, \omega))^2 dt \rightarrow 0 \text{ as } n \rightarrow +\infty,$$

for each realization  $\omega$  fixed. Hence, by bounded convergence theorem we also have

$$\mathbb{E} \left( \int_S^T (h - \psi_n)^2 ds \right) \rightarrow 0 \text{ as } n \rightarrow +\infty,$$

**Step 3.** Finally, given any  $f \in V$  we find a sequence  $h_n \in V$  such that each  $h_n(t, \omega)$  is bounded, and

$$\mathbb{E} \left( \int_S^T (f - h_n)^2 ds \right) \rightarrow 0 \text{ as } n \rightarrow +\infty.$$

This is done as follows: set

$$h_n = \begin{cases} -n & \text{if } f(t, \omega) \leq -n, \\ f(t, \omega) & \text{if } -n \leq f(t, \omega) \leq n, \\ n & \text{if } f(t, \omega) \geq n. \end{cases}$$

Note that on the set  $\{f \geq n\}$  we have

$$(f - h_n)^2 = (f - n)^2 \leq f^2,$$

and similarly for the set  $\{f \leq -n\}$ . It follows that

$$\mathbb{E} \left( \int_S^T (f - h_n)^2 ds \right) \leq \mathbb{E} \left( \int_S^T f^2(t, \omega) \chi_{|f(t, \omega)| \geq n} dt \right) \rightarrow 0 \text{ as } n \rightarrow +\infty,$$

by the Lebesgue dominated convergence theorem, since

$$\mathbb{E} \left( \int_S^T f^2(t, \omega) dt \right) < +\infty.$$

Together, steps (1)-(3) show that for any function  $f \in V$  we can find a sequence of bounded elementary functions  $\phi_n$  such that

$$\mathbb{E} \left( \int_S^T (f - \phi_n)^2 ds \right) \rightarrow 0 \text{ as } n \rightarrow +\infty, \quad (2.11)$$

that is,  $\phi_n$  converges to  $f$  in the space  $L^2([S, T] \times \Omega)$ . Then, given  $f \in V$  we choose such sequence  $\phi_n$  and define

$$\int_S^T f(t, \omega) dB_t(\omega) = \lim_{n \rightarrow +\infty} \int_S^T \phi_n(t, \omega) dB_t. \quad (2.12)$$

We now need to check two things: (i) the limit in the right side exists for any sequence  $\phi_n$  for which (2.11) holds, and (ii) it does not depend on the particular choice of the sequence  $\phi_n$ . Actually, (ii) follows from (i) (this is a simple exercise). The reason the limit exists is that the sequence

$$\alpha_n(\omega) = \int_S^T \phi_n(t, \omega) dB_t$$

is a Cauchy sequence in  $L^2(\Omega)$ . Indeed, we have

$$\mathbb{E}(\alpha_n - \alpha_m)^2 = \mathbb{E} \left( \int_S^T (\alpha_n(t, \omega) - \alpha_m(t, \omega))^2 dt \right) \rightarrow 0 \text{ as } n, m \rightarrow +\infty, \quad (2.13)$$

because the sequence  $\phi_n(t)$  is convergent (hence, Cauchy) in  $L^2([S, T] \times \Omega)$ . The first equality in (2.13) follows from Ito's isometry.

**Corollary 2.4** *If  $f(t, \omega) \in V$  then*

$$\mathbb{E} \left( \int_S^T f(t, \omega) dB_t \right)^2 = \mathbb{E} \left( \int_S^T f^2(t, \omega) dt \right).$$

**Corollary 2.5** *If  $f(t, \omega) \in V$  and  $f_n(t, \omega) \in V$  are such that*

$$\mathbb{E} \left( \int_S^T (f_n(t, \omega) - f(t, \omega))^2 dt \right) \rightarrow 0,$$

*then*

$$\int_S^T f_n(t, \omega) dB_t \rightarrow \int_S^T f(t, \omega) dB_t \text{ in } L^2(\Omega).$$

### An explicit example

Let us show that

$$\int_0^t B_s dB_s = \frac{1}{2} B_t^2 - \frac{1}{2} t. \quad (2.14)$$

Consider the elementary functions

$$\phi_n(s, \omega) = \sum_j B(t_j, \omega) \chi_{[t_j, t_{j+1})}(s),$$

where  $t_j = jt/n$ ,  $j = 0, \dots, n-1$ . Then, we have

$$\begin{aligned} \mathbb{E} \left( \int_0^t (\phi_n(s) - B_s)^2 ds \right) &= \mathbb{E} \sum_j \int_{t_j}^{t_{j+1}} (B(t_j) - B(s))^2 ds = \sum_j \int_{t_j}^{t_{j+1}} (t_j - s) ds \\ &= \sum_j \frac{(t_{j+1} - t_j)^2}{2} \rightarrow 0, \end{aligned}$$

as  $\Delta t_j \rightarrow 0$ . Hence,

$$\int_0^t B_s dB_s = \lim_{n \rightarrow +\infty} \int_0^t \phi_n(s, \omega) dB_s = \lim_{n \rightarrow +\infty} \sum_{j=0}^{n-1} B(t_j) (B(t_{j+1}) - B(t_j)),$$

with the limit understood in  $L^2(\Omega)$ -sense. Note that

$$B^2(t_{j+1}) - B^2(t_j) = (B(t_{j+1}) - B(t_j))^2 + 2B(t_j)(B(t_{j+1}) - B(t_j)),$$

and thus

$$\sum_{j=0}^{n-1} B(t_j) (B(t_{j+1}) - B(t_j)) = \frac{1}{2} \sum_{j=0}^{n-1} (B^2(t_{j+1}) - B^2(t_j)) - \frac{1}{2} \sum_{j=0}^{n-1} (B(t_{j+1}) - B(t_j))^2. \quad (2.15)$$

The first sum above is telescoping:

$$\sum_{j=0}^{n-1} (B^2(t_{j+1}) - B^2(t_j)) = B^2(t).$$

For the second sum in the right side of (2.15) we claim that

$$\sum_{j=0}^{n-1} (B(t_{j+1}) - B(t_j))^2 \rightarrow t \text{ as } j \rightarrow +\infty, \text{ in } L^2(\Omega). \quad (2.16)$$

This is verified by a direct computation:

$$\mathbb{E} \left( \sum_j (B(t_{j+1}) - B(t_j))^2 - t \right)^2 = \mathbb{E} \sum_{j,k} ((B(t_{j+1}) - B(t_j))^2 - \Delta t_j) ((B(t_{k+1}) - B(t_k))^2 - \Delta t_k). \quad (2.17)$$

Independence of the increments of the Brownian motion implies that the terms with  $j \neq k$  in (2.17) vanish since

$$\mathbb{E}((B(t_{k+1}) - B(t_k))^2 - \Delta t_k) = 0,$$

and the same for  $k$  replaced by  $j$ . It follows that

$$\begin{aligned} \mathbb{E} \left( \sum_j (B(t_{j+1}) - B(t_j))^2 - t \right)^2 &= \mathbb{E} \left( \sum_j ((B(t_{j+1}) - B(t_j))^2 - \Delta t_j)^2 \right) \\ &= \sum_j \mathbb{E} ((B(t_{j+1}) - B(t_j))^4 + (t_{j+1} - t_j)^2 - 2(B(t_{j+1}) - B(t_j))^2(t_{j+1} - t_j)) \\ &= \sum_j (3(t_{j+1} - t_j)^2 + (t_{j+1} - t_j)^2 - 2(t_{j+1} - t_j)^2) = 2 \sum_j (t_{j+1} - t_j)^2 \rightarrow 0 \text{ as } \Delta t_j \rightarrow 0. \end{aligned} \quad (2.18)$$

Therefore, (2.16) holds, and we have proved (2.14).

## Martingales

**Definition 2.6** Let  $\mathcal{F}_t$  be a family of  $\sigma$ -algebras such that  $\mathcal{F}_s \subseteq \mathcal{F}_t$  for all  $0 \leq s \leq t$ . A random process  $M_t$  is an  $\mathcal{F}_t$ -martingale if

- (i) The process  $M_t$  is  $\mathcal{F}_t$ -measurable for all  $t \geq 0$ .
- (ii) The expectation  $\mathbb{E}(|M_t|) < +\infty$  for all  $t \geq 0$ , and
- (iii) The conditional expectation  $\mathbb{E}(M_s | \mathcal{F}_t) = M_t$  a.s. for all  $s \geq t$ .

The main non-technical assumption here is the last one: the conditional expectation of a martingale is its present value. One may think of a martingale as a ‘‘fair game.’’ If  $M_t$  represents a gambler’s account balance at time  $t$ , then the condition  $E[M_t | \mathcal{F}_s] = M_s$  says that the expected future balance, given the current balance, is unchanged. So the game favors neither the gambler nor the house. Of course, the change  $M_t - M_s$  may be positive or negative, but its expected value conditioned on  $M_s$  is zero. On the other hand, if  $M_t$  is a sub-martingale, which means that (iii) above is replaced by

$$\mathbb{E}(M_s | \mathcal{F}_t) \geq M_t \text{ for all } s \geq t, \quad (2.19)$$

then given the gambler’s current account balance, he may expect his earnings to increase. If  $X_t$  is a super-martingale, that is,

$$\mathbb{E}(M_s | \mathcal{F}_t) \leq M_t \text{ for all } s \geq t, \quad (2.20)$$

then given the gambler's current account balance, he may expect his earnings to decrease (this seems to be the most realistic model, given the success of many casinos).

A useful observation is that if  $X_s$  is an  $\mathcal{F}_s$ -martingale and a function  $\phi(x)$  is convex, then  $\phi(X_s)$  is a sub-martingale. This follows from Jensen's inequality. For instance, if  $X_s$  is a martingale,  $|X_s|$ ,  $X_s^2$ , and all  $X_s^{2m}$  with  $m \in \mathbb{N}$ , are sub-martingales. Continuous martingales satisfy a remarkable property that estimates the maximum of a process by the terminal time statistics.

**Theorem 2.7** (*Continuous Doob inequality*) *If  $M_t$  is a continuous in time martingale such that  $\mathbb{E}(|M_T|^p) < +\infty$ , then for all  $p \geq 1$ ,  $T \geq 0$  and  $\lambda > 0$  we have*

$$P \left[ \sup_{0 \leq t \leq T} |M_t| \geq \lambda \right] \leq \frac{1}{\lambda^p} \mathbb{E}(|M_T|^p).$$

We will not prove this result here but rather prove it only for discrete martingales. A sequence  $X_j$  is a martingale with respect to a sequence of  $\sigma$ -algebras  $\mathcal{F}_j$  if (i)  $\mathcal{F}_n \subseteq \mathcal{F}_{n+1}$ , (ii)  $X_n$  is  $\mathcal{F}_n$ -measurable, (iii)  $E[|X_n|] < +\infty$ , and (iv)  $E(X_{n+1}|\mathcal{F}_n) = X_n$  almost surely. It follows that  $E(X_m|\mathcal{F}_n) = X_n$  almost surely for all  $m \geq n$ . The discrete Doob's inequality is the following estimate that bounds the supremum of  $X_j$  in terms of the expectation of the last element:

**Theorem 2.8** (*Discrete Doob's inequality*) *Suppose  $(X_j, \mathcal{F}_j)$ ,  $1 \leq j \leq n$ , is a martingale sequence such that  $\mathbb{E}(|X_n|^p) < +\infty$ , then for any  $l > 0$  and any  $p \geq 1$  we have*

$$P \left\{ \omega : \sup_{1 \leq j \leq n} |X_j| \geq l \right\} \leq \frac{1}{l^p} E(|X_n|^p).$$

**Proof.** Let us define  $S(\omega) = \sup_{1 \leq j \leq n} |X_j(\omega)|$ . Then the event  $E = \{\omega : S(\omega) \geq l\}$  can be decomposed as a disjoint union of the sets

$$E_j = \{\omega : |X_1(\omega)| < l, \dots, |X_{j-1}(\omega)| < l, |X_j(\omega)| \geq l\},$$

that is,  $E = \bigcup_{j=1}^n E_j$  and  $E_j \cap E_m = \emptyset$  for  $j \neq m$ . Note that, as  $|X_j| \geq l$  on the set  $E_j$  we have an inequality

$$P(E_j) \leq \frac{1}{l^p} \int_{E_j} |X_j|^p dP.$$

The function  $\phi(x) = |x|^p$  is convex for  $p \geq 1$ , hence, as we mentioned above, the sequence  $|X_j|^p$  is a sub-martingale, thus  $|X_j|^p \leq \mathbb{E}(|X_n|^p|\mathcal{F}_j)$ , and

$$P(E_j) \leq \frac{1}{l^p} \int_{E_j} |X_j|^p dP \leq \frac{1}{l^p} \int_{E_j} \mathbb{E}(|X_n|^p|\mathcal{F}_j) dP$$

Moreover, the set  $E_j$  is  $\mathcal{F}_j$ -measurable as follows immediately from the way  $E_j$  is defined, hence

$$P(E_j) \leq \frac{1}{l^p} \int_{E_j} \mathbb{E}(|X_n|^p|\mathcal{F}_j) dP = \frac{1}{l^p} \int_{E_j} |X_n|^p dP,$$

simply from the definition of the conditional expectation  $\mathbb{E}(|X_n|^p|\mathcal{F}_j)$ . Now, summing over all  $j$  and using the fact that  $E_j$  are disjoint we obtain

$$P(E) = \sum_{j=1}^n P(E_j) \leq \frac{1}{l^p} \sum_{j=1}^n \int_{E_j} |X_n|^p dP \leq \frac{1}{l^p} \int_E |X_n|^p dP \leq \frac{1}{l^p} \int_{\Omega} |X_n|^p dP = \frac{1}{l^p} \mathbb{E}(|X_n|^p),$$

and we are done.  $\square$

## Ito integral as a martingale

It turns out that Ito integral is always a martingale (this is a great advantage of the Ito integral compared to Stratonovich and other definitions of the stochastic integral).

**Theorem 2.9** *Let  $f \in V(0, T)$  for all  $T > 0$  then the process  $M_t(\omega) = \int_0^t f(s, \omega) dB_s(\omega)$  is an  $\mathcal{F}_t$ -martingale, and*

$$P \left( \sup_{0 \leq t \leq T} |M_t| \geq \lambda \right) \leq \frac{1}{\lambda^2} \mathbb{E} \left( \int_0^T f^2(s, \omega) ds \right),$$

for all  $\lambda > 0$ .

Another important property of the Ito integral is that it has a continuous in time version.

**Theorem 2.10** *Let  $f \in V(0, T)$  then there exists a  $t$ -continuous version of*

$$M_t(\omega) = \int_0^t f(s, \omega) dB_s(\omega).$$

We will not prove these results here (the proofs are not long or difficult – see Oksendal's book) – both are consequences of the continuous Doob inequality.

In order to understand why Theorems 2.9 and 2.10 hold let us consider an elementary function

$$f(t) = \sum_j e_j(\omega) \chi_{[t_j, t_{j+1})}(t).$$

Then, for each  $t$  we find  $N$  such that  $t \in [t_N, t_{N+1}]$  and write

$$M_t = \int_0^t f(s, \omega) dB_s(\omega) = \sum_{j=1}^{N-1} e_j(\omega) (B_{t_{j+1}} - B_{t_j}) + e_N(\omega) (B_t - B_{t_N}). \quad (2.21)$$

Since  $B_t$  is almost surely continuous, it follows immediately that  $M_t(\omega)$  is almost surely continuous on each interval  $(t_j, t_{j+1})$ . It is also clear that there is no jump at the times  $t_j$ , hence  $M_t(\omega)$  is a.s. continuous for all  $t \geq 0$ . In order to verify that

$$\mathbb{E}(M_t | \mathcal{F}_s) = M_s, \quad (2.22)$$

let us assume that  $s = 0$ , then (2.22) is simply

$$\mathbb{E}(M_t) = M_0 = 0. \quad (2.23)$$

Using (2.24) we write

$$\begin{aligned} \mathbb{E}(M_t) &= \mathbb{E} \left( \int_0^t f(s, \omega) dB_s(\omega) \right) = \sum_{j=1}^{N-1} \mathbb{E}(e_j(\omega) (B_{t_{j+1}} - B_{t_j})) + \mathbb{E}(e_N(\omega) (B_t - B_{t_N})) \\ &= \sum_{j=1}^{N-1} \mathbb{E}(e_j(\omega)) \mathbb{E}(B_{t_{j+1}} - B_{t_j}) + \mathbb{E}(e_N(\omega)) \mathbb{E}(B_t - B_{t_N}) = 0, \end{aligned} \quad (2.24)$$

hence (2.23), indeed, holds. Once again, we used above the independence of  $e_j(\omega)$  and the forward increment  $B_{t_{j+1}} - B_{t_j}$  (recall that this is true because  $f(t, \omega)$  is  $\mathcal{F}_t$ -adapted). If  $s > 0$  we use essentially the same argument starting from

$$M_t - M_s = \int_s^t f(\tau, \omega) dB_\tau. \quad (2.25)$$

## The Ito formula

We begin with the definition of the Ito processes.

**Definition 2.11** *A one-dimensional Ito process is a stochastic process of the form*

$$X_t = X_0 + \int_0^t u(s, \omega) ds + \int_0^t v(s, \omega) dB_s,$$

with  $u(s, \omega)$  and  $v(s, \omega)$  such that

$$P \left( \int_0^t v^2(s, \omega) ds < +\infty \text{ for all } t \geq 0. \right) = 1,$$

and

$$P \left( \int_0^t |u(s, \omega)| ds < +\infty \text{ for all } t \geq 0. \right) = 1.$$

A shorter notation is

$$dX_t = u dt + v dB_t.$$

**Example.** Our computation of  $\int_0^t B_s dB_s$  shows that the process  $X_t = B_t^2$  can be written as

$$B_t^2 = t + 2 \int_0^t B_s dB_s,$$

or

$$dX_t = dt + 2B_t dB_t.$$

The Ito formula gives a recipe on how to express a function of an Ito process as another Ito process.

**Theorem 2.12** *(The Ito formula) Let  $X_t$  be an Ito process given by*

$$dX_t = u dt + v dB_t,$$

and let  $g(t, x) \in C^2([0, +\infty) \times \mathbb{R})$ . Then  $Y_t = g(t, X_t)$  is also an Ito process, and

$$dY_t = \frac{\partial g(t, X_t)}{\partial t} dt + \frac{\partial g(t, X_t)}{\partial x} dX_t + \frac{1}{2} \frac{\partial^2 g(t, X_t)}{\partial x^2} (dX_t)^2.$$

Here  $(dX_t)^2 = dX_t \cdot dX_t$  with the convention  $dt \cdot dt = dt \cdot dB_t = dB_t \cdot dt = 0$ , and  $dB_t \cdot dB_t = dt$ .

### An interlude on the random walk

The Ito formula has an extra term involving  $\partial^2 g(t, X_t) / \partial x^2$  in the right side that is absent in the deterministic case  $v = 0$ . As this term is what really matters in the connection to elliptic and parabolic partial differential equations, let us explain where it comes from. We need to look no further than the random walk that approximates the Brownian motion. Recall that the correct scaling is to take the time step  $\Delta t = h^2$  and the spatial step  $\Delta x = h$ , with  $h$  small. That is, the random walk satisfies

$$P(X_{t_{n+1}} = X_{t_n} \pm h | X_{t_n}) = 1/2,$$

where  $t = nh^2$  and  $x = mh$  with some integers  $n$  and  $m$ . Then, given a twice continuously differentiable function  $f(t, x)$ , and since  $t_{n+1} - t_n = h^2$ , we have

$$\begin{aligned} f(t_{n+1}, X_{t_{n+1}}) &= f(t_n, X_{t_n}) + \frac{\partial f(t_n, X_{t_n})}{\partial x}(X_{t_{n+1}} - X_{t_n}) \\ &+ \frac{\partial f(t_n, X_{t_n})}{\partial t}h^2 + \frac{1}{2} \frac{\partial^2 f(t_n, X_{t_n})}{\partial x^2}(X_{t_{n+1}} - X_{t_n})^2 \\ &+ \frac{\partial^2 f(t_n, X_{t_n})}{\partial t^2}h^4 + \frac{\partial^2 f(t_n, X_{t_n})}{\partial x \partial t}h^2(X_{t_{n+1}} - X_{t_n}) + \dots \end{aligned} \quad (2.26)$$

Note that  $(X_{t_{n+1}} - X_{t_n})^2 = h^2$  with probability one – this is where the convention  $dB_t \cdot dB_t = dt$  comes from! Hence the two terms in the second line above are of the same order, and we can actually rewrite it as

$$\begin{aligned} f(t_{n+1}, X_{t_{n+1}}) &= f(t_n, X_{t_n}) + \frac{\partial f(t_n, X_{t_n})}{\partial x}(X_{t_{n+1}} - X_{t_n}) \\ &+ \left[ \frac{\partial f(t_n, X_{t_n})}{\partial t} + \frac{1}{2} \frac{\partial^2 f(t_n, X_{t_n})}{\partial x^2} \right] (t_{n+1} - t_n) \\ &+ \frac{\partial^2 f(t_n, X_{t_n})}{\partial t^2}(t_{n+1} - t_n)^2 + \frac{\partial^2 f(t_n, X_{t_n})}{\partial x \partial t}(t_{n+1} - t_n)(X_{t_{n+1}} - X_{t_n}) + \dots \end{aligned} \quad (2.27)$$

The terms in the last line in (2.27) vanish as  $h \rightarrow 0$ . Summing over  $n$  and formally passing to the limit  $h \rightarrow 0$  in (2.27) we get

$$f(B_t) = f(B_0) + \int_0^t \frac{\partial f(s, B_s)}{\partial x} dB_s + \int_0^t \left[ \frac{\partial f(t, B_t)}{\partial t} + \frac{1}{2} \frac{\partial^2 f(t, B_t)}{\partial x^2} \right] ds.$$

This is, of course, nothing but the Ito formula.

### Sketch of proof of Ito's formula

Let us assume that

$$dX_t = udt + vdB_t.$$

Using Taylor's formula we get, for any partition  $0 = t_0 < t_1 < \dots < t_N = t$ , and  $X_j = X_{t_j}$

$$\begin{aligned} g(t, X_t) &= g(0, X_0) + \sum_j (g(t_{j+1}, X_{j+1}) - g(t_j, X_j)) \\ &= g(0, X_0) + \sum_j \frac{\partial g(t_j, X_j)}{\partial t} \Delta t_j + \sum_j \frac{\partial g(t_j, X_j)}{\partial x} \Delta X_j \\ &+ \frac{1}{2} \sum_j \frac{\partial^2 g(t_j, X_j)}{\partial t^2} (\Delta t_j)^2 + \sum_j \frac{\partial^2 g(t_j, X_j)}{\partial t \partial x} \Delta t_j \Delta X_j + \frac{1}{2} \sum_j \frac{\partial^2 g(t_j, X_j)}{\partial x^2} (\Delta X_j)^2 + \sum_j R_j, \end{aligned}$$

where  $R_j = o(|\Delta t_j|^2 + |\Delta X_j|^2)$ . If we let  $\Delta t_j \rightarrow 0$ , we get

$$\sum_j \frac{\partial g(t_j, X_j)}{\partial t} \Delta t_j \rightarrow \int_0^t \frac{\partial g(s, X_s)}{\partial s} ds,$$

and

$$\sum_j \frac{\partial g(t_j, X_j)}{\partial x} \Delta X_j \rightarrow \int_0^t \frac{\partial g(s, X_s)}{\partial x} dX_s.$$

In order to understand what happens to the term

$$\frac{1}{2} \sum_j \frac{\partial^2 g(t_j, X_j)}{\partial x^2} (\Delta X_j)^2,$$

we assume that  $u$  and  $v$  are elementary. In the general case we may approximate  $u$  and  $v$  by elementary functions and use a density argument. If  $u$  and  $v$  are elementary, then, after possibly refining the partition  $\{t_j\}$  to make sure that  $u_j$  and  $v_j$  are constant on the intervals  $(t_j, t_{j+1})$ , we have

$$\Delta X_j = u_j \Delta t_j + v_j \Delta B_j,$$

hence

$$\begin{aligned} \sum_j \frac{\partial^2 g(t_j, X_j)}{\partial x^2} (\Delta X_j)^2 &= \sum_j \frac{\partial^2 g(t_j, X_j)}{\partial x^2} u_j^2 (\Delta t_j)^2 + 2 \sum_j \frac{\partial^2 g(t_j, X_j)}{\partial x^2} u_j v_j (\Delta t_j) \Delta B_j \\ &+ \sum_j \frac{\partial^2 g(t_j, X_j)}{\partial x^2} v_j^2 (\Delta B_j)^2 = I + II + III. \end{aligned}$$

It is easy to check that

$$\mathbb{E}(I^2) \rightarrow 0, \quad \mathbb{E}(II^2) \rightarrow 0 \text{ as } \Delta t_j \rightarrow 0.$$

Indeed, we have:

$$\mathbb{E}(II^2) = 4 \sum_{i,j} \mathbb{E} \left( \frac{\partial^2 g(t_j, X_j)}{\partial x^2} \frac{\partial^2 g(t_i, X_i)}{\partial x^2} u_j v_j u_i v_i (\Delta t_j) \Delta B_j (\Delta t_i) \Delta B_i \right) = 8 \sum_{i>j} + 4 \sum_{i=j}.$$

But using our beloved independence of the increments from the past and the fact that  $u_j$  and  $v_j$  are  $\mathcal{F}_t$ -adapted, we get that

$$\begin{aligned} &\sum_{i>j} \mathbb{E} \left( \frac{\partial^2 g(t_j, X_j)}{\partial x^2} \frac{\partial^2 g(t_i, X_i)}{\partial x^2} u_j v_j u_i v_i (\Delta t_j) \Delta B_j (\Delta t_i) \Delta B_i \right) \\ &= \sum_{i>j} \mathbb{E} \left( \frac{\partial^2 g(t_j, X_j)}{\partial x^2} \frac{\partial^2 g(t_i, X_i)}{\partial x^2} u_j v_j u_i v_i (\Delta t_j) \Delta B_j (\Delta t_i) \right) \mathbb{E}(\Delta B_i) = 0. \end{aligned}$$

Therefore, we have, again using independence of the increments from the past and the fact that  $u_j$  and  $v_j$  are  $\mathcal{F}_t$ -adapted:

$$\begin{aligned} \mathbb{E}(II^2) &= 4 \sum_j \mathbb{E} \left( \left( \frac{\partial^2 g(t_j, X_j)}{\partial x^2} \right)^2 u_j^2 v_j^2 (\Delta t_j)^2 (\Delta B_j)^2 \right) \\ &= 4 \sum_j \mathbb{E} \left( \left( \frac{\partial^2 g(t_j, X_j)}{\partial x^2} \right)^2 u_j^2 v_j^2 (\Delta t_j)^2 \right) \mathbb{E}((\Delta B_j)^2) \\ &= 4 \sum_j \mathbb{E} \left( \left( \frac{\partial^2 g(t_j, X_j)}{\partial x^2} \right)^2 u_j^2 v_j^2 \right) (\Delta t_j)^3 \rightarrow 0 \text{ as } \Delta t_j \rightarrow 0. \end{aligned}$$

A very similar computation shows that  $\mathbb{E}(I^2) \rightarrow 0$  as  $\Delta t_j \rightarrow 0$ . The last step is to show that

$$III \rightarrow \int_0^t \frac{\partial^2 g(s, X_s)}{\partial x^2} v^2(s, \omega) ds \text{ in } L^2(\Omega) \text{ as } \Delta t_j \rightarrow 0. \quad (2.28)$$

In order to check this, we set

$$a(t) = \frac{\partial^2 g(t, x)}{\partial x^2} v^2(t, \omega),$$

and  $a_j = a(t_j)$ . Consider then

$$\mathbb{E}\left(\sum_j a_j (\Delta B_j)^2 - \sum_j a_j \Delta t_j\right)^2 = \sum_{ij} \mathbb{E}(a_i a_j ((\Delta B_i)^2 - \Delta t_i) ((\Delta B_j)^2 - \Delta t_j)). \quad (2.29)$$

As before, if  $i > j$  then the forward increment  $\Delta B_i$  is independent of the other terms in (2.29). Since  $\mathbb{E}(\Delta B_i^2) = \Delta t_i$ , the terms with  $i \neq j$  in (2.29) vanish, and we get

$$\mathbb{E}\left(\sum_j a_j (\Delta B_j)^2 - \sum_j a_j \Delta t_j\right)^2 = \sum_j \mathbb{E}\left(a_j^2 ((\Delta B_j)^2 - \Delta t_j)^2\right). \quad (2.30)$$

As  $a_j$  is  $\mathcal{F}_{t_j}$  measurable, we deduce that  $\Delta B_j$  is independent from  $a_j$ , hence

$$\mathbb{E}\left(\sum_j a_j (\Delta B_j)^2 - \sum_j a_j \Delta t_j\right)^2 = \sum_j \mathbb{E}(a_j^2) \mathbb{E}((\Delta B_j)^2 - \Delta t_j)^2 \quad (2.31)$$

$$\begin{aligned} &= \sum_j \mathbb{E}(a_j^2) \mathbb{E}((\Delta B_j)^4 - 2(\Delta B_j)^2 \Delta t_j + (\Delta t_j)^2) = \sum_j \mathbb{E}(a_j^2) (3(\Delta t_j)^2 - 2(\Delta t_j)^2 + (\Delta t_j)^2) \\ &= 2 \sum_j \mathbb{E}(a_j^2) (\Delta t_j)^2 \rightarrow 0 \text{ as } \Delta t_j \rightarrow 0. \end{aligned} \quad (2.32)$$

**Example 1:** Let  $X_t = B_t$  and  $\phi(x) = x^2$ . Then by applying Itô's formula to the process  $Y_t = (B_t)^2$  we find that

$$(B_t)^2 = \int_0^t 2B_s dB_s + t,$$

as we have seen already.

**Example 2:** Let  $X_t = B_t$  and  $\phi(x) = e^{\alpha x}$ . Itô's formula applied to the process  $Y_t = e^{\alpha B_t}$  gives

$$Y_t = 1 + \int_0^t \alpha Y_s dB_s + \frac{\alpha^2}{2} \int_0^t Y_s ds$$

which may be expressed as

$$dY_t = \alpha Y_t dB_t + \frac{\alpha^2}{2} Y_t dt$$

Similarly, the function  $Z_t = e^{\alpha B_t - \alpha^2 t/2}$  satisfies

$$dZ_t = \alpha Z_t dB_t$$

This shows that  $Z_t$  is a martingale since

$$Z_t = 1 + \int_0^t \alpha Z_s dB_s$$

and the Itô integral is a martingale. (Actually one can easily compute directly that  $Z$  is a martingale without the help of stochastic calculus.) It follows that

$$\mathbb{E}(Z_t) = Z_0 = 1,$$

whence

$$\mathbb{E}(e^{\alpha B_t}) = e^{\alpha^2 t/2}. \quad (2.33)$$

There are many other ways to compute the expectation in (2.33) but this probably is the simplest.

**Theorem 2.13 (Itô Product Rule)** *Suppose that  $X_t(\omega)$  and  $Y_t(\omega)$  two stochastic processes satisfying*

$$\begin{aligned} dX_t &= F(X_t, t)dt + G(X_t, t)dB_t \\ dY_t &= H(Y_t, t)dt + K(Y_t, t)dB_t. \end{aligned}$$

Then the process  $Z_t(\omega) = X_t(\omega)Y_t(\omega)$  satisfies

$$\begin{aligned} dZ_t &= (F(X_t, t)Y_t + H(Y_t, t)X_t + G(X_t, t)K(Y_t, t)) dt + (G(X_t, t)Y_t + K(Y_t, t)X_t)dB_t \\ &= Y_t dX_t + X_t dY_t + G(X_t, t)K(Y_t, t)dt. \end{aligned}$$

### Itô's formula in multiple dimensions

We can also define vector-valued stochastic integrals using a  $m$ -dimensional Brownian motion. Suppose that  $G(s, \omega)$  is a matrix valued process such that

$$G^{ij}(s, \omega) \in \mathcal{L}^2([0, T]), \quad i = 1, \dots, d, \quad j = 1, \dots, m.$$

If  $B_t$  is a  $m$ -dimensional Brownian motion, then

$$X_t = \int_0^t G(s, \omega) dB_t$$

defines a  $d$ -dimensional stochastic process whose components are

$$X_t^{(i)} = \sum_{j=1}^m \int_0^t G^{ij}(s, \omega) dB_t^{(j)}, \quad i = 1, \dots, d$$

Itô's formula extends to multiple dimensions in the following way.

**Theorem 2.14** Suppose that  $B_t$  is a  $m$ -dimensional Brownian motion and that  $X_t(\omega) = (X_t^{(i)}(\omega))_i$  is a  $d$ -dimensional stochastic process satisfying

$$X_t^{(i)}(\omega) = X_0^{(i)}(\omega) + \int_0^t F^{(i)}(s, \omega) ds + \sum_{j=1}^m \int_0^t G^{ij}(s, \omega) dB_t^{(j)},$$

If  $\phi(x_1, \dots, x_d, t)$  is twice-differentiable in the spatial variables, differentiable in  $t$ , then the one-dimensional process  $Y_t := \phi(X_t(\omega), t)$  satisfies

$$\begin{aligned} dY_t &= [F(t, \omega) \cdot \nabla \phi(X_t(\omega), t) + \phi_t(X_t(\omega), t)] dt + \sum_{j=1}^m \sum_{i=1}^d \frac{\partial \phi}{\partial x_i}(X_t(\omega), t) G^{ij}(t, \omega) dB_t^{(j)} \\ &+ \frac{1}{2} \left( \sum_{k=1}^m \sum_{i,j=1}^d \phi_{x_i x_j}(X_t(\omega), t) G^{(ik)}(t, \omega) G^{(jk)}(t, \omega) \right) dt. \end{aligned}$$

You can remember the last term by Taylor's formula and the heuristic formula

$$dB_t^{(i)} dB_t^{(j)} \sim dt, \quad \text{if } i = j, \quad (= 0, \text{ otherwise})$$

so that

$$(G^{(kh)} dB_t^{(h)})(G^{(qp)} dB_t^{(p)}) \sim \delta_{hp} G^{(kh)} G^{(qp)} dt$$

Thus, off-diagonal terms ( $p \neq h$ ) vanish in the formula.

## Stochastic differential equations

**Example.** Let us consider a stochastic differential equation

$$dN_t = rN_t dt + \alpha N_t dB_t. \tag{2.34}$$

Ito's formula says that

$$d(\ln N_t) = \frac{dN_t}{N_t} - \frac{1}{2N_t^2} (dN_t)^2 = \frac{dN_t}{N_t} - \frac{\alpha^2}{2N_t^2} N_t^2 dt = \frac{dN_t}{N_t} - \frac{\alpha^2}{2} dt.$$

It follows that

$$\int_0^t \frac{dN_s}{N_s} = \ln N_t - \ln N_0 + \alpha^2 t.$$

On the other hand, (2.34) implies that

$$\int_0^t \frac{dN_s}{N_s} = rt + \alpha B_t.$$

Therefore, we have an explicit solution

$$N_t = N_0 \exp\left(\left(r - \frac{\alpha^2}{2}\right)t + \alpha B_t\right).$$

As a consequence, we also have

$$\mathbb{E}(N_t) = N_0 e^{rt},$$

as can also be seen immediately from (2.34). An interesting property of  $N_t$  is that if  $r < \alpha^2/2$  then  $N_t \rightarrow 0$  as  $t \rightarrow \infty$ , almost surely, even though  $\mathbb{E}(N_t) \rightarrow +\infty$ .

In the general case we have the following result.

**Theorem 2.15** Let  $T > 0$ , and  $b(t, x)$  and  $\sigma(t, x)$  satisfy

$$|b(t, x)| + |\sigma(t, x)| \leq C(1 + |x|),$$

and

$$|b(t, x) - b(t, y)| + |\sigma(t, x) - \sigma(t, y)| \leq D|x - y|,$$

for all  $x, y \in \mathbb{R}^n$ . Then the stochastic differential equation

$$dX_t = b(t, X_t)dt + \sigma(t, X_t)dB_t, \quad X_0 = x_0,$$

has a unique continuous in  $t$  solution  $X_t(\omega)$  that is  $\mathcal{F}_t$ -adapted, and

$$\mathbb{E} \left( \int_0^T |X_t|^2 dt \right) < +\infty.$$

The Lipschitz continuity and at most linear growth of  $b$  and  $\sigma$  are needed even for existence and uniqueness theorems for ordinary differential equations. For instance, solutions of the ODE  $\dot{X} = X^2$  blow up in a finite time if  $X(0) > 0$  while the ODE  $\dot{X} = 2\sqrt{X}$  with the initial data  $X(0) = 0$  has two solutions:  $X_1(t) \equiv 0$ , and  $X_2(t) = t^2$ .

### 3 Representations of solutions of PDEs

We now develop representation formulas for solutions of various PDEs – we will for now take existence and regularity of solutions for granted but will later address them separately.

#### Poisson's equation in the whole space

Consider an elliptic operator

$$\mathcal{L}f(x) = \frac{1}{2} \sum_{i,j=1}^n a_{ij}(x) \frac{\partial^2 f}{\partial x_i \partial x_j} + \sum_{i=1}^n b_i(x) \frac{\partial f}{\partial x_i}. \quad (3.1)$$

We assume that  $a_{ij}(x)$  are sufficiently smooth and bounded, and the matrix  $a_{ij}(x)$  is symmetric:  $a_{ij} = a_{ji}$ . A more important assumption is that the operator  $\mathcal{L}$  is uniformly elliptic. This means that there exists a constant  $c > 0$  so that for all  $\xi \in \mathbb{R}^n$  we have

$$\sum_{i,j=1}^n a_{ij}(x) \xi_i \xi_j \geq c|\xi|^2. \quad (3.2)$$

Let  $\sigma_{ij}$  be a matrix such that  $a = \sigma\sigma^T$ , and each component of  $\sigma$  is bounded and in  $C^1(\mathbb{R}^n)$ . Consider  $X_t$ , the solution to the SDE

$$X_t = x + \int_0^t b(X_s)ds + \int_0^t \sigma(X_s)dB_s.$$

The Ito formula for a function  $f(X_t)$  is

$$f(X_t) - f(X_0) = \int_0^t \mathcal{L}f(s, X_s)ds + \sum_{j,k=1}^n \int_0^t \sigma_{jk}(s, X_s) \frac{\partial f(s, X_s)}{\partial x_j} dB_k(s). \quad (3.3)$$

**Theorem 3.1** Let  $\lambda > 0$  and  $f(x)$  be a  $C^1$  function of compact support. Suppose  $u(x)$  is a  $C_b^2(\mathbb{R}^n)$  solution of the Poisson equation

$$-\mathcal{L}u + \lambda u(x) = f(x), \quad x \in \mathbb{R}^n. \quad (3.4)$$

Then

$$u(x) = \mathbb{E}_x \int_0^\infty e^{-\lambda t} f(X_t) dt. \quad (3.5)$$

**Proof.** Let  $u(x)$  be the solution to (3.4), then, using the Ito formula (3.3) we obtain

$$u(X_t) - u(X_0) = \int_0^t \mathcal{L}u(X_s) ds + \sum_{j,k=1}^n \int_0^t \sigma_{jk}(s, X_s) \frac{\partial u(s, X_s)}{\partial x_j} dB_k(s).$$

Moreover, if we set  $v(t, x) = e^{-\lambda t} u(x)$ , then

$$\begin{aligned} v(t, X_t) - v(0, X_0) &= \int_0^t e^{-\lambda s} \mathcal{L}u(X_s) ds - \lambda \int_0^t e^{-\lambda s} u(X_s) ds \\ &\quad + \sum_{j,k=1}^n \int_0^t e^{-\lambda s} \sigma_{jk}(s, X_s) \frac{\partial u(s, X_s)}{\partial x_j} dB_k(s). \end{aligned}$$

Taking the expectation gives

$$e^{-\lambda t} \mathbb{E}(u(X_t)) - u(x) = - \int_0^t e^{-\lambda s} \mathbb{E}(f(X_s)) ds.$$

Letting  $t \rightarrow +\infty$  leads to (3.5) since the functions  $u(x)$  and  $f(x)$  bounded.  $\square$

In order to see why we need to take  $\lambda > 0$  in (3.4), consider the one-dimensional case:

$$-u'' + k^2 u = f(x), \quad (3.6)$$

so that  $\mathcal{L} = d^2/dx^2$ , and  $X_t$  is simply the one-dimensional Brownian motion  $B_t$ . If  $k = 0$  then (3.6) need not have a bounded solution (though it might for some  $f(x)$ ) – this is easily seen if  $f(x) \geq 0$  since then  $u(x)$  is a concave function and that contradicts the fact that  $u(x)$  is bounded unless  $u(x) \equiv \text{const}$ . This can also be seen from the probabilistic formula (3.5): if  $\lambda = 0$  and  $f(x) \geq 0$  is a compactly supported function then the integral in (3.5) diverges since  $B_t$  is recurrent in one dimension. On the other hand, when  $k > 0$ , a bounded solution of (3.6) is given by an explicit formula

$$u(x) = \frac{1}{2k} \int_{-\infty}^{\infty} e^{-k|x-y|} f(y) dy. \quad (3.7)$$

An interesting exercise is to compare (3.5) and (3.7) (with  $\lambda = k^2$ ) and deduce the law of the Brownian motion in one dimension.

## The Feynman-Kac Formula

Let us now look at time dependent problems. In what follows, we will work with solutions to initial value problems and with solutions to terminal value problems. One can switch between these two perspectives through a simple change of variables:  $t \rightarrow T - t$ . Suppose that  $w(t, x) \in C^{2,1}([0, \infty) \times \mathbb{R})$  solves the initial value problem

$$w_t = \frac{\sigma^2(x)}{2}w_{xx} + b(x)w_x \quad x \in \mathbb{R}, \quad t > 0 \quad (3.8)$$

with initial data  $w(x, 0) = f(x)$ , which is smooth and compactly supported. We also assume that  $b(x)$  and  $\sigma(x)$  are Lipschitz continuous and bounded, and that  $w$  is bounded. Then, for  $t > 0$  fixed, the function  $u(s, x) = w(t - s, x)$  satisfies the terminal value problem

$$u_s + \frac{\sigma^2(x)}{2}u_{xx} + b(x)u_x = 0, \quad x \in \mathbb{R}, \quad s < t \quad (3.9)$$

with terminal condition  $u(t, x) = f(x)$ . Moreover,  $u \in C^{2,1}((-\infty, t] \times \mathbb{R})$ . Now, let  $B_s(\omega)$  be a standard Brownian motion with filtration  $(\mathcal{F}_s)_{s \geq 0}$ . Suppose that,  $X_s(\omega)$  is an  $\mathcal{F}_s$ -adapted solution to the stochastic ODE

$$dX_s = b(X_s)ds + \sigma(X_s)dB_s \quad (3.10)$$

with the initial condition  $X_0 = x$ . The existence and uniqueness of such a solution is guaranteed by our assumptions about  $b$  and  $\sigma$ .

Now, a direct application of Ito's formula shows us that,

$$\begin{aligned} u(t, X_t) - u(0, X_0) &= \int_0^t \left( u_s(s, X_s) + b(X_s)u_x(s, X_s) + \frac{\sigma^2(X_s)}{2}u_{xx}(s, X_s) \right) ds \\ &+ \int_0^t \sigma(X_s)u_x(s, X_s)dB_s = \int_0^t \sigma(X_s)u_x(s, X_s)dB_s. \end{aligned} \quad (3.11)$$

We used the PDE (3.9) in the last step. Therefore, taking the expectation, we find that

$$\mathbb{E}[u(t, X_t)] = \mathbb{E}[u(0, X_0)] = u(0, x), \quad (3.12)$$

since the Itô integral has zero mean. In terms of  $w$ , this shows that

$$w(t, x) = u(0, x) = \mathbb{E}[u(t, X_t)] = \mathbb{E}[f(X_t)] \quad (3.13)$$

In summary, these arguments demonstrate the following:

**Theorem 3.2** (i) **Initial value problem:** *Suppose that  $w(t, x) \in C^{2,1}([0, \infty) \times \mathbb{R})$  is bounded and satisfies*

$$w_t = \frac{\sigma^2(x)}{2}w_{xx} + b(x)w_x \quad x \in \mathbb{R}, \quad t > 0 \quad (3.14)$$

*with initial condition  $w(0, x) = f(x) \in C_0^2(\mathbb{R})$ . Then  $w(t, x)$  is represented by*

$$w(t, x) = \mathbb{E}_x[f(X_t)] \quad (3.15)$$

where

$$dX_s = b(X_s)ds + \sigma(X_s)dB_s \text{ for } s \geq 0 \text{ and } X_0(\omega) = x.$$

(ii) **Terminal value problem:** Suppose that  $u(t, x) \in C^{2,1}((-\infty, T] \times \mathbb{R})$  is bounded and satisfies

$$u_t + \frac{\sigma^2(x)}{2}u_{xx} + b(x)u_x = 0 \quad x \in \mathbb{R}, \quad t < T \quad (3.16)$$

with terminal condition  $u(T, x) = f(x) \in C_0^2(\mathbb{R})$ . Then  $u(0, x)$  is represented by

$$u(0, x) = E[f(X_T)]. \quad (3.17)$$

If we need to find  $u(t, x)$  for  $t \in (0, T)$  we simply consider the process

$$dX_s = b(X_s)ds + \sigma(X_s)dB_s,$$

that starts at time  $s = t$  at the point  $x$ :  $X_t = x$ . Then we have

$$u(t, x) = \mathbb{E}(f(X_T)).$$

## Generalizations

To avoid technical difficulties, we have been rather conservative in our assumptions about the initial conditions and the coefficients. In fact, these representations hold under milder conditions on the initial data and the coefficients. Now let us suppose that  $u(t, x)$  satisfies the second-order linear PDE

$$u_t + \sum_{i,j=1}^d \frac{1}{2}a_{ij}(t, x)u_{x_i x_j} + \sum_{j=1}^d b_j(t, x)u_{x_j} + c(t, x)u = 0, \quad x \in \mathbb{R}^d, \quad t < T \quad (3.18)$$

with terminal condition  $u(x, T) = f(x)$  which is continuous (but not necessarily differentiable or bounded). We also assume

- The matrix  $a_{ij}$  is given by  $a_{ij} = \sum_k \sigma_{ik}\sigma_{kj} = \sigma\sigma^T$  for some matrix  $\sigma_{jk}(t, x)$ .
- The matrix  $a_{ij} = a_{ij}(t, x)$  is uniformly positive definite:  $\sum_{ij} a_{ij}\xi_j\xi_i \geq \mu|\xi|^2$  for some constant  $\mu > 0$ , independent of  $(t, x)$ .
- Both  $\sigma_{ij}(t, x)$  and  $b(t, x) = (b_j(t, x))$  are Lipschitz continuous in  $x$ , continuous in  $t$ , and grow at most linearly in  $x$ .
- The function  $c(t, x)$  is continuous in  $(t, x)$  and bounded in  $x$ .
- The terminal condition  $f(x)$  satisfies the growth condition  $|f(x)| \leq Ce^{p|x|^2}$  for some constant  $p > 0$  sufficiently small.
- $u(t, x)$  satisfies the growth condition  $|u(t, x)| \leq Ce^{p|x|^2}$  for  $x \in \mathbb{R}$ ,  $t \in [t_0, T]$  and some constant  $p > 0$  sufficiently small.

Suppose that for a given  $(x, t)$ , the process  $X_s^{x,t}(\omega) : [t, T] \times \Omega \rightarrow \mathbb{R}^d$  satisfies

$$dX_s^{x,t} = b(s, X_s^{x,t}) ds + \sum_j \sigma_{ij}(s, X_s^{x,t}) dB_s^{(j)}, \quad s \in [t, T], \quad (3.19)$$

with  $X_t^{x,t}(\omega) = x$ . The superscripts indicate that the process  $X_s^{x,t}$  starts at the point  $x$  at time  $t$ . Notice that  $X_s^{x,t}$  is a vector, and  $b(s, x)$  is also a vector. Then one can prove:

**Theorem 3.3** *Under the assumptions given above,  $u(x, t)$  satisfies*

$$u(t, x) = \mathbb{E} \left[ f(X_T^{x,t}) e^{\int_t^T c(X_s^{x,t}, s) ds} \right]. \quad (3.20)$$

**Sketch of proof:** To prove this statement, one may apply Itô's formula and the product rule to the process defined by

$$H_r(\omega) = u(X_r^{x,t}, r) e^{\int_t^r c(X_s^{x,t}, s) ds}, \quad r \in [t, T]. \quad (3.21)$$

The fact that the terminal data  $f(x)$  may not be smooth or bounded causes some difficulty that may be overcome by using Itô's formula with stopping times. For  $n > 0$ , let  $S_n(\omega)$  be the stopping time  $S_n = \inf\{s \geq t \mid |X_s^{x,t}| \geq n\}$ . Then we conclude that for  $r \in (t, T)$ ,

$$\begin{aligned} H_{r \wedge S_n} - H_t &= u(X_{r \wedge S_n}^{x,t}, r \wedge S_n) e^{\int_t^{r \wedge S_n} c(X_s^{x,t}, s) ds} - u(X_t^{x,t}, t) \\ &= \int_t^{r \wedge S_n} e^{\int_t^s c(X_\tau^{x,t}, \tau) d\tau} \left( u_s + \sum_j b_j u_{x_j} + \frac{1}{2} a_{ij} u_{x_i x_j} + c(X_s^{x,t}, s) u \right) ds \\ &\quad + \int_t^{r \wedge S_n} e^{\int_t^s c(X_\tau^{x,t}, \tau) d\tau} \sum_{i,j} u_{x_i} \sigma_{ij} dB_s^{(j)} \\ &= \int_t^{r \wedge S_n} e^{\int_t^s c(X_\tau^{x,t}, \tau) d\tau} \sum_{i,j} u_{x_i} \sigma_{ij} dB_s^{(j)} \quad (\text{using (3.18)}) \end{aligned}$$

Notice that arguments inside the integrals are evaluated at  $(X_s^{x,t}, s)$ . Taking the expectation as before, we conclude that

$$u(x, t) = \mathbb{E} [u(X_t^{x,t}, t)] = \mathbb{E} \left[ u(X_{r \wedge S_n}^{x,t}, r) e^{\int_t^{r \wedge S_n} c(X_s^{x,t}, s) ds} \right]. \quad (3.22)$$

Notice that if  $u$  itself is not bounded, then the expectation on the right is not obviously finite. This explains our use of the stopping time – the stopping time restricts  $X_{r \wedge S_n}^{x,t}$  to a bounded region, over which  $u$  must be bounded since  $u$  is continuous. The next step is to take  $n \rightarrow \infty$ . Using the growth assumptions on  $u$  and the coefficients one can show that as  $n \rightarrow \infty$ , the above expression remains finite since  $P(S_n < r) = O(e^{-\alpha n^2})$  as  $n \rightarrow \infty$ . This shows that

$$u(x, t) = \mathbb{E} \left[ u(X_r^{x,t}, r) e^{\int_t^r c(X_s^{x,t}, s) ds} \right] \quad (3.23)$$

Then we let  $r \rightarrow T$ . If we knew that  $u$  were sufficiently smooth and bounded at  $r = T$ , then we could apply Itô's formula with  $r = T$  in the above formula. This was our approach in the

first section, since we assumed the initial (or terminal) data was  $C^2$ . In general, however, this is not the case. Nevertheless, one may use the dominated convergence theorem to show that as  $r \rightarrow T$ ,

$$\lim_{r \rightarrow T} \mathbb{E} \left[ u(X_r^x, r) e^{\int_t^r c(X_s^x, s) ds} \right] = \mathbb{E} \left[ f(X_T^x) e^{\int_t^T c(X_s^x, s) ds} \right] \quad (3.24)$$

even when  $f$  is merely continuous and satisfies a growth condition (see Karatzas and Shreve for more details).  $\square$

Next we formulate a similar result for the initial value problem. Suppose that  $w(x, t)$  satisfies

$$w_t = \sum_{i,j=1}^d \frac{1}{2} a_{ij}(x, t) u_{x_i x_j} + \sum_{j=1}^d b_j(x, t) u_{x_j} + c(x, t) u, \quad x \in \mathbb{R}^d, \quad t > 0 \quad (3.25)$$

with initial condition  $w(x, 0) = f(x)$ . Then the function  $\tilde{u}(x, -t) := w(x, t)$  satisfies (3.19) with  $T = 0$ , and coefficients given by  $\tilde{a}_{ij}(x, t) = a_{ij}(x, -t)$ ,  $\tilde{b}(x, t) = b(x, -t)$ ,  $\tilde{c}(x, t) = c(x, -t)$ . For given  $(x, t)$  let  $X_s^{x,t}(\omega)$  satisfy

$$dX_s^{x,t} = b(X_s^{x,t}, t-s) ds + \sum_j \sigma_{ij}(X_s^{x,t}, t-s) dB_s^{(j)}, \quad s \in [0, t] \quad (3.26)$$

Then the analysis above shows that

$$w(x, t) = \mathbb{E} \left[ f(X_t^{x,t}) e^{\int_0^t c(X_s^{x,t}, t-s) ds} \right]. \quad (3.27)$$

In particular, if  $c \equiv 0$ , then

$$w(x, t) = E \left[ f(X_t^{x,t}) \right]. \quad (3.28)$$

These are very elegant formulas which have a natural physical interpretation. Here is how one can think about it. The equation (3.25) models the diffusion, transport, and reaction of a scalar quantity  $w(x, t)$ . The vector field  $b$  is the “drift” or wind. The matrix  $a_{ij}$  determines the rates of diffusion in a given direction. The process  $X_t^{x,t}$  may be thought of as the paths of particles diffusing in this velocity field. The function  $c(x, t)$  represents a reaction rate. So, imagine hot, reactive particles being carried in the wind. Now, consider the formula (3.28) for the case  $c \equiv 0$  (no reaction). What determines the temperature at a point  $(x, t)$ ? The temperature at this point is determined by which particles arrive at point  $x$  at time  $t$  and how hot those particles were initially. The quantity  $f(X_t^{x,t})$  represents the initial “temperature” evaluated at the “end” of the path  $f(X_t^{x,t})$ . Notice that  $X_s^{x,t}$  actually runs backwards in the time-frame associated with the PDE. Roughly speaking,  $f(X_t^{x,t})$  tells us what information propagates to the point  $x$  at time  $t$ . The paths are random; formula (3.28) says that the solution is determined by the expectation over all such particles. In the case that  $c \neq 0$ , formula (3.27) tells us that the reaction heats up each particle along its trajectory, increasing (or decreasing) its temperature by a factor of  $e^{\int_0^t c(X_s^{x,t}, t-s) ds}$ . Notice that when  $a_{ij}$ ,  $b$ , and  $c$  are independent of  $t$ , we can replace  $t-s$  in the above expressions with  $s$ .

## Poisson's equation

Here we use Itô's formula to derive a representation for solutions to Poisson's equation with a variable zero-order term. Suppose that  $w(x)$  is  $C^2$  and bounded, and satisfies

$$\sum_{i,j} \frac{1}{2} a_{ij}(x) w_{x_i x_j} + \sum_j b_j(x) w_{x_j} - c(x) w = f(x), \quad x \in \mathbb{R}^d \quad (3.29)$$

with  $c(x) \geq c_0 > 0$  for some constant  $c_0 > 0$ . As before, we assume  $a_{ij} = \sigma \sigma^T$  is uniformly positive definite, and that  $a$ ,  $b$ , and  $c$  satisfy the continuity criteria given earlier.

**Theorem 3.4** *Suppose that  $X_t(\omega)$  solves the stochastic differential equation*

$$dX_t = b(X_t) dt + \sum_j \sigma_{ij}(X_t) dB_t^{(j)}, \quad t \geq 0 \quad (3.30)$$

with  $X_0(\omega) = x \in \mathbb{R}^d$ , almost surely. The solution  $w(x)$  is represented by

$$w(x) = \mathbb{E} \left[ \int_0^\infty e^{-\int_0^s c(X_\tau) d\tau} f(X_s) ds \right]. \quad (3.31)$$

**Proof:** Now apply Itô's formula and the product rule to the process

$$H_t(\omega) = e^{-\int_0^t c(X_s) ds} w(X_t^x). \quad (3.32)$$

We compute:

$$\begin{aligned} H_t - H_0 &= w(X_t) e^{-\int_0^t c(X_s) ds} - w(X_0) \\ &= \int_0^t e^{-\int_0^s c(X_\tau) d\tau} \left( \sum_j b_j w_{x_j} + \frac{1}{2} \sum_{i,j} a_{ij} w_{x_i x_j} - c(X_s) w \right) ds \\ &\quad + \int_0^t e^{-\int_0^s c(X_\tau) d\tau} \sum_{i,j} w_{x_i} \sigma_{ij} dB_s^{(j)} \\ &= \int_0^t e^{-\int_0^s c(X_\tau) d\tau} f(X_s) ds + \int_0^t e^{-\int_0^s c(X_\tau) d\tau} \sum_{i,j} w_{x_i} \sigma_{ij} dB_s^{(j)}. \end{aligned} \quad (3.33)$$

Now we take the expectation of both sides and let  $t \rightarrow \infty$ . Due to the lower bound on  $c(x)$ ,

$$\lim_{t \rightarrow \infty} \left| \mathbb{E} \left[ w(X_t) e^{-\int_0^t c(X_s) ds} \right] \right| \leq \lim_{t \rightarrow \infty} e^{-c_0 t} \|w\|_\infty = 0. \quad (3.34)$$

Therefore,

$$w(x) = -\mathbb{E} \int_0^\infty e^{-\int_0^s c(X_\tau) d\tau} f(X_s^x) ds, \quad (3.35)$$

and we are done.  $\square$

## Problems in bounded domains

So far we have considered solutions to partial differential equations posed in the whole space  $x \in \mathbb{R}^d$ . Itô's formula also leads to representation formulas for solutions to PDE's posed in a bounded domain with appropriate boundary conditions. We consider two types of problems: boundary value problems for elliptic equations and initial value/terminal value problems for parabolic equations.

### Boundary value problems for elliptic operators

Suppose that  $D \subset \mathbb{R}^d$  is a smooth, bounded domain. Let  $w(x) \in C^2(\bar{D})$  satisfy

$$\sum_{i,j} \frac{1}{2} a_{ij}(x) w_{x_i x_j} + \sum_j b_j(x) w_{x_j} - c(x) w = f(x), \quad x \in D, \quad (3.36)$$

with boundary condition  $w(x) = g(x)$  for  $x \in \partial D$ . The function  $g(x)$  is prescribed. As usual, we assume that the operator in the left side is elliptic and coefficients are sufficiently regular and bounded.

In addition, we need to assume that the function  $c(x) \geq 0$ . The need for this extra assumption can be seen on the very simple one-dimensional example: consider the problem

$$u'' - cu = 0, \quad u(0) = 0, \quad u(\pi) = 0.$$

If  $c = -1$  this problem has two linearly independent solutions:  $u_1(x) \equiv 0$ , and  $u_2(x) = \sin x$ . On the other hand, if  $c(x) \geq 0$  this can not happen: solution of the boundary value problem is always unique. Indeed, let  $w_1$  and  $w_2$  be two solutions of (3.36) with some prescribed functions  $f(x)$  and  $g(x)$ . The difference  $w(x) = w_1(x) - w_2(x)$  satisfies

$$\sum_{i,j} \frac{1}{2} a_{ij}(x) w_{x_i x_j} + \sum_j b_j(x) w_{x_j} - c(x) w = 0, \quad x \in D, \quad (3.37)$$

with the boundary condition  $w(x) = 0$  for  $x \in \partial D$ . Since  $c(x) \geq 0$  the maximum principle applies and shows that  $w(x) = 0$  in  $D$  meaning that solution is unique. The reason why the maximum principle applies if  $c(x) \geq 0$  is easily seen if we impose a slightly stronger condition  $c(x) > 0$ . Then, if  $w(x)$  solves (3.37) and attains its maximum at some interior point  $x_0 \in D$ , the Hessian matrix  $H(x)$  with the entries

$$H_{ij}(x_0) = \frac{\partial^2 w(x_0)}{\partial x_i \partial x_j}$$

is non-positive definite. Therefore, since  $a_{ij}(x_0)$  is a positive-definite matrix (that follows from ellipticity), we have

$$\sum_{ij} a_{ij}(x_0) \frac{\partial^2 w(x_0)}{\partial x_i \partial x_j} \leq 0.$$

In addition, at  $x_0$  we have

$$\frac{\partial w(x_0)}{\partial x_j} = 0$$

for all  $j = 1, \dots, n$ . Using this in (3.37) we obtain

$$-c(x)w(x_0) \geq 0,$$

whence  $w(x_0) \leq 0$  if  $c(x) > 0$  in  $D$ . It follows that  $w(x) \leq 0$  for all  $x \in D$ . A similar analysis at the point  $x_1$  where  $w(x)$  attains its minimum shows that

$$-c(x)w(x_1) \leq 0,$$

meaning that  $w(x_1) \geq 0$ . Therefore,  $w(x) \geq 0$  in  $D$ , and we conclude that  $w(x) = 0$ . The weaker assumption  $c(x) \geq 0$  requires a slightly more subtle analysis but the basic idea is the same.

How can we represent the solution of the boundary value problem? If  $X_t(\omega)$  solves the stochastic differential equation

$$dX_t = b(X_t) dt + \sum_j \sigma_{ij}(X_t) dB_t^{(j)}, \quad t \geq 0 \quad (3.38)$$

with  $X_0(\omega) = x \in D$ , then the trajectories will travel outside of the set  $D$ , where the function  $w$  is not defined. To overcome this difficulty, we define the stopping time

$$\gamma_D(\omega) = \inf\{t \mid X_t \in \mathbb{R}^d \setminus D\}.$$

This is the first hitting time to the boundary  $\partial D$ . A basic result in the SDE theory (see Richard Bass book "Diffusions and Elliptic Operators", Proposition I.8.2) says that  $\gamma_D(\omega) < +\infty$  a.s. if the domain  $D$  is bounded. Then, we can define the process

$$H_t(\omega) = e^{-\int_0^{t \wedge \gamma_D} c(X_s) ds} w(X_{t \wedge \gamma_D}^x). \quad (3.39)$$

Here, we denote  $t \wedge \gamma_D := \min(t, \gamma_D)$ . Itô's formula and the product rule then imply that

$$\begin{aligned} H_t - H_0 &= w(X_{t \wedge \gamma_D}) e^{-\int_0^{t \wedge \gamma_D} c(X_s) ds} - w(X_0) \\ &= \int_0^{t \wedge \gamma_D} e^{-\int_0^s c(X_\tau) d\tau} \left( \sum_j b_j w_{x_j} + \frac{1}{2} \sum_{i,j} a_{ij} w_{x_i x_j} - c(X_s) w \right) ds \\ &\quad + \int_0^{t \wedge \gamma_D} e^{-\int_0^s c(X_\tau) d\tau} \sum_{i,j} w_{x_i} \sigma_{ij} dB_s^{(j)} \\ &= \int_0^{t \wedge \gamma_D} e^{-\int_0^s c(X_\tau) d\tau} f(X_s) ds + \int_0^{t \wedge \gamma_D} e^{-\int_0^s c(X_\tau) d\tau} \sum_{i,j} w_{x_i} \sigma_{ij} dB_s^{(j)} \end{aligned} \quad (3.40)$$

As before, we now take the expectation of both sides and let  $t \rightarrow \infty$ . As  $\gamma_D(t)$  is finite a.s., we have  $\lim_{t \rightarrow \infty} \gamma_D(\omega) \wedge t = \gamma_D(\omega)$ , also almost surely. Consequently, the fact that  $w$  is bounded and that  $c \geq 0$ , allows us to use the dominated convergence theorem to show that

$$\lim_{t \rightarrow \infty} \mathbb{E} \left[ w(X_{t \wedge \gamma_D}) e^{-\int_0^{t \wedge \gamma_D} c(X_s) ds} \right] = \mathbb{E} \left[ w(X_{\gamma_D}) e^{-\int_0^{\gamma_D} c(X_s) ds} \right] = \mathbb{E} \left[ g(X_{\gamma_D}) e^{-\int_0^{\gamma_D} c(X_s) ds} \right]. \quad (3.41)$$

Similarly, using the fact that  $f$  is bounded and  $E[\gamma_D] < \infty$ , we may use the dominated convergence theorem to show that

$$\lim_{t \rightarrow \infty} \mathbb{E} \left[ \int_0^{t \wedge \gamma_D} e^{-\int_0^s c(X_\tau) d\tau} f(X_s) ds \right] = \mathbb{E} \left[ \int_0^{\gamma_D} e^{-\int_0^s c(X_\tau) d\tau} f(X_s) ds \right] \quad (3.42)$$

Therefore, taking  $t \rightarrow \infty$ , we obtain a representation for  $w(x)$ :

$$w(x) = E \left[ g(X_{\gamma_D}) e^{-\int_0^{\gamma_D} c(X_s) ds} \right] - E \left[ \int_0^{\gamma_D} e^{-\int_0^s c(X_\tau) d\tau} f(X_s) ds \right]. \quad (3.43)$$

Notice that with the stronger assumption  $c(x) \geq c_0 > 0$ , we could lift the condition that  $E[\gamma_D] < \infty$ , which was used in the application of the dominated convergence theorem to obtain (3.42). We could also lift the restriction that  $w \in C^2(\bar{D})$ , and require only that  $w \in C^2(D) \cap C(\bar{D})$  (thus, the second derivatives might blow up at that boundary). To handle this case, stop the process when it is distance  $\epsilon$  from the boundary. Then let  $\epsilon \rightarrow 0$ .

**Example 1:** In particular, this representation shows that if  $w(x)$  solves  $\Delta w = 0$  in  $D$  with  $w(x) = g(x)$  for  $x \in \partial D$ , then

$$w(x) = E \left[ g(x + \sqrt{2}B_{\gamma_D}) \right] \quad (3.44)$$

The quantity  $g(x + \sqrt{2}B_{\gamma_D})$  is the boundary function evaluated at the point where the process first hits the boundary. The solution to the PDE is the expectation of these values.

## Initial boundary value problems

Suppose that  $D \subset \mathbb{R}^d$  is a smooth bounded domain. Let  $D_T = D \times (0, T]$  denote the parabolic cylinder. Suppose that  $w(t, x) \in C^{2,1}(D_T) \cap C(\bar{D}_T)$  satisfies the initial value problem

$$\begin{aligned} w_t &= \sum_{i,j} \frac{1}{2} a_{ij}(t, x) w_{x_i x_j} + \sum_j b_j(t, x) w_{x_j} + c(t, x) w, \quad x \in D, \quad t > 0 \\ w(0, x) &= f(x) \quad x \in D \\ w(t, x) &= g(t, x) \quad x \in \partial D, \quad t \geq 0. \end{aligned}$$

Here we assume that  $c(x, t)$  is bounded and continuous. For given  $(t, x) \in D_T$ , let  $X_s^{x,t}(\omega)$  satisfy

$$dX_s^{x,t} = b(t-s, X_s^{x,t}) ds + \sum_j \sigma_{ij}(t-s, X_s^{x,t}) dB_s^{(j)}, \quad s \in [0, t]. \quad (3.45)$$

Define the stopping time  $\gamma_D^{x,t} = \inf\{s \geq 0 \mid X_s^{x,t} \in \mathbb{R} \setminus D\}$ . This is the first time the process hits the boundary of the set  $D$ . Then define  $\gamma^{x,t} = \gamma_D^{x,t} \wedge t$ . This is also a stopping time, and it represents the time at which the process  $(X_s^{x,t}, t-s)$  hits the parabolic boundary  $(D \times \{0\}) \cup (\partial D \times [0, T])$ , which is the boundary of the set  $D_T$ . For convenient notation, let us define the function

$$k(t, x) = \begin{cases} f(x), & \text{if } t = 0, x \in \bar{D} \\ g(t, x), & \text{if } t > 0, x \in \partial D \end{cases} \quad (3.46)$$

This function is equal to  $f(x)$  at the base of the parabolic boundary, and it is equal to  $g(t, x)$  on the sides of the parabolic boundary.

**Theorem 3.5** *Under the above assumptions,  $w(x, t)$  satisfies*

$$w(x, t) = E \left[ k(X_{\gamma^{x,t}}, \gamma^{x,t}) e^{\int_0^{\gamma^{x,t}} c(X_s^{x,t}, t-s) ds} \right] \quad (3.47)$$

**Proof:** I leave this as an exercise. It may be proved as in the other cases.  $\square$

## Transition Densities

Consider the vector-valued stochastic process defined by

$$dX_t = b(X_t) dt + \sigma^{ij}(X_t) dW_t^j \quad \text{for } t > 0, \quad X_0(\omega) = x. \quad (3.48)$$

Suppose that  $a_{ij} = \sigma \sigma^T$  is uniformly positive. Suppose also that  $a$  and  $b$  satisfy the continuity conditions described previously. Because of the Markov property of Brownian motion, one can show that  $X_t$  is a Markov process satisfying

$$P(X_t \in A | \mathcal{F}_s) = P(X_t \in A | X_s), \quad \forall s \in [0, t]. \quad (3.49)$$

Suppose that  $X_t$  has a smooth transition density  $p(x, s; y, t)$ . This means that

$$P(X_t \in A | X_s = x) = \int_A p(x, s; y, t) dy \quad (3.50)$$

and

$$E[f(X_t) | X_s = x] = \int_{\mathbb{R}^d} f(y) p(x, s; y, t) dy \quad (3.51)$$

for suitable functions  $f$ . What equation does  $p(x, s; y, t)$  satisfy?

Here is a formal computation that can be made rigorous under suitable smoothness and growth assumptions on the coefficients  $b$  and  $\sigma^{ij}$ . If  $f(x)$  is smooth and compactly supported, then Itô's formula tells us that

$$f(X_t) - f(X_s) = \int_s^t \mathcal{A}f(X_r) dr + \int_s^t \sum_{ij} \frac{\partial f}{\partial x_i}(X_r) \sigma^{ij}(X_r) dW_r^j \quad (3.52)$$

where  $\mathcal{A}$  denotes the differential operator

$$\mathcal{A}f(y) := \frac{1}{2} \sum_{ij} a_{ij}(y) f_{y_i y_j} + b(y) \cdot \nabla_y f \quad (3.53)$$

Conditioning on the event  $X_s = x$  and taking the expectation, we obtain

$$E[f(X_t) | X_s = x] - E[f(X_s) | X_s = x] = \int_s^t E[\mathcal{A}f(X_r) | X_s = x] dr. \quad (3.54)$$

Now using the definition of the transition density, we may write this expression as

$$\int_{\mathbb{R}^d} f(y)p(x, s; y, t) dy - f(x) = \int_s^t \int_{\mathbb{R}^d} (\mathcal{A}f(y))p(x, s; y; r) dy dr. \quad (3.55)$$

Formally differentiating both sides with respect to  $t$ , we obtain the equation

$$\int_{\mathbb{R}^d} f(y)p_t(x, s; y, t) dy = \int_{\mathbb{R}^d} (\mathcal{A}f(y))p(x, s; y; t) dy. \quad (3.56)$$

Now, integrate by parts on the right hand side:

$$\begin{aligned} \int_{\mathbb{R}^d} (\mathcal{A}f(y))p(x, s; y; t) dy &= \int_{\mathbb{R}^d} \left( \frac{1}{2} \sum_{ij} a_{ij}(x) f_{y_i y_j} + b(x) \cdot \nabla_y f \right) p(x, s; y; t) dy \\ &= \int_{\mathbb{R}^d} f(y) \left( \frac{1}{2} \sum_{ij} \frac{\partial^2}{\partial y_i \partial y_j} (a_{ij}(x)p(x, s; y; t)) - \nabla_y \cdot (b(y)p(x, s; y; t)) \right) dy \\ &= \int_{\mathbb{R}^d} f(y) (\mathcal{A}_y^* p(x, s; y, t)) dy. \end{aligned} \quad (3.57)$$

Here  $\mathcal{A}_y^*$  is the **adjoint operator** defined by

$$\mathcal{A}_y^* g(y) := \frac{1}{2} \sum_{ij} \frac{\partial^2}{\partial y_i \partial y_j} (a_{ij}(y)g(y)) - \nabla_y \cdot (b(y)g(y)). \quad (3.58)$$

In the integration by parts step, the boundary terms vanish since  $f$  has compact support. Therefore,  $p(x, s; y, t)$  should satisfy

$$\int_{\mathbb{R}^d} f(y) (p_t(x, s; y, t) - \mathcal{A}_y^* p(x, s; y, t)) dy = 0. \quad (3.59)$$

Since  $f(y)$  is chosen arbitrarily, and since we assume  $p$  to be sufficiently smooth, this implies that for each fixed  $x$  and  $s$ , the function  $u(y, t) = p(x, s; y, t)$  satisfies  $u_t = \mathcal{A}_y^* u$ . That is,

$$\frac{\partial}{\partial t} p(x, s; y, t) = \mathcal{A}_y^* p(x, s; y, t). \quad (3.60)$$

As  $t \searrow s$ ,  $p(x, s; y, t)$  as a function of  $y$  converges to a delta distribution centered at  $y = x$ . This equation (3.60) is often called the **Kolmogorov forward equation** for the transition density  $p(x, s; y, t)$ . The term “forward” is applied since it describes the forward evolution of the probability density for  $X_t$ .

For fixed  $y$  and  $t$ , the function  $u(x, s) = p(x, s; y, t)$  satisfies a different equation. To derive this equation, suppose that  $f$  is again smooth and compactly supported. We have already shown that the solution to the terminal value problem

$$w_s + \mathcal{A}_x w = 0, \quad s < t, \quad x \in \mathbb{R}^d \quad (3.61)$$

with terminal data  $w(x, t) = f(x)$  has the representation

$$w(x, s) = E[f(X_t) | X_s = x] = \int_{\mathbb{R}^d} f(y)p(x, s; y, t) dy \quad (3.62)$$

Formally differentiating the integral expression with respect to  $s$  and  $x$  and using (3.61) we find that

$$\int_{\mathbb{R}^d} f(y)(p_s(x, s; y, t) dy + \mathcal{A}_x p(x, s; y, t)) dy = 0 \quad (3.63)$$

Since  $f$  was arbitrarily chosen this implies that for each  $y$ ,  $p_s(x, s; y, t) + \mathcal{A}_x p(x, s; y, t) = 0$ . Since  $x$  and  $s$  were also arbitrarily chosen, this suggests that for each  $y$  and  $t$  fixed,

$$\frac{\partial}{\partial s} p(x, s; y, t) + \mathcal{A}_x p(x, s; y, t) = 0. \quad (3.64)$$

Since the coefficients defining the process  $X_t$  are independent of  $t$ , the transition density is a function of  $t - s$ :

$$p(x, s; y, t) = \rho(x, y, t - s) \quad (3.65)$$

for some function  $\rho(x, y, r)$ . Then (3.64) shows that for fixed  $y$ ,  $\rho(x, y, t)$  satisfies

$$\frac{\partial}{\partial t} \rho(x, y, t) = \mathcal{A}_x \rho(x, y, t) \quad (3.66)$$

This equation is often called the **Kolmogorov backward equation**.

## 4 Second order elliptic equations

### 4.1 Sobolev spaces

#### Weak derivatives

A weak derivative is a natural extension of the derivative to a non-differentiable function. In order to motivate this notion, let  $u \in C^1(\mathbb{R}^n)$  and  $\phi$  be a smooth compactly supported test function. Then we have:

$$\int_{\mathbb{R}^n} u \frac{\partial \phi}{\partial x_i} = - \int_{\mathbb{R}^n} \frac{\partial u}{\partial x_i} \phi dx.$$

Note that the left side makes sense whether  $u$  is differentiable or not – all we need is that  $u(x)$  is in  $L^1_{loc}(\mathbb{R}^n)$ , that is,  $u$  is integrable over compact sets. This motivates the following definition.

**Definition 4.1** *Let  $u, v \in L^1_{loc}(\mathbb{R}^n)$ , then  $v = \partial u / \partial x_j$  is the weak derivative of  $u$  with respect to  $x_j$  if for any  $\phi \in C_c^\infty(\mathbb{R}^n)$  we have*

$$\int_{\mathbb{R}^n} u \frac{\partial \phi}{\partial x_i} dx = - \int_{\mathbb{R}^n} v \phi dx.$$

**Example 1.** Let  $u(x) = |x|$ , then for any  $\phi \in C_c^\infty(\mathbb{R})$  we have

$$\int_{\mathbb{R}} |x| \phi'(x) dx = \int_{-\infty}^0 (-x) \phi'(x) dx + \int_0^{\infty} x \phi'(x) dx = \int_{-\infty}^0 \phi(x) dx - \int_0^{\infty} \phi(x) dx = \int_{\mathbb{R}} \text{sgn}(x) \phi(x) dx.$$

Therefore, the weak derivative  $u'(x) = \text{sgn}x$ .

**Example 2.** Let  $u(x) = \text{sgn}(x)$ , we claim that the weak derivative of  $u(x)$  does not exist. Indeed, for any test function  $\phi \in C_c^\infty(\mathbb{R})$  we have

$$\int_{\mathbb{R}} u\phi' dx = \int_0^\infty \phi'(x) dx - \int_{-\infty}^0 \phi'(x) dx = -2\phi(0).$$

Let us now assume that there exists some function  $v \in L_{loc}^1(\mathbb{R})$  such that for any function  $\phi$  as above we have

$$\int_{\mathbb{R}} v(x)\phi(x) dx = 2\phi(0). \quad (4.1)$$

Choose a sequence  $\phi_m(x)$  such that  $\phi_m(0) = 1$ ,  $\phi_m(x) = 0$  for  $|x| \geq 1/m$  and  $0 \leq \phi_m(x) \leq 1$  for all  $x \in [-1/m, 1/m]$ . Then we have

$$\left| \int_{\mathbb{R}} v(x)\phi_m(x) dx \right| \leq \int_{-1/m}^{1/m} |v(x)| dx \rightarrow 0,$$

as  $m \rightarrow +\infty$ , since  $v \in L_{loc}^1(\mathbb{R})$ . This contradicts (4.1) and shows that the weak derivative  $u'(x)$  does not exist.

The definition of the higher order weak derivatives is a natural extension of the above. If  $\alpha = (\alpha_1, \dots, \alpha_m)$  is a multi-index and  $|\alpha| = \alpha_1 + \dots + \alpha_m$ , then  $v = D^\alpha u$  in a domain  $U$  if  $v \in L_{loc}^1(U)$ , and for any function  $\phi \in C_c^\infty(U)$  we have

$$\int_U u D^\alpha \phi dx = (-1)^{|\alpha|} \int_U v \phi dx.$$

**Definition 4.2** *The Sobolev space  $W^{k,p}(U)$  consists of all functions  $u \in L_{loc}^1(U)$  such that for each multi-index  $\alpha$  with  $|\alpha| \leq k$ , the weak derivative  $D^\alpha u$  exists and belongs to  $L^p(U)$ .*

When  $p = 2$  we use a special notation  $H^k(U) = W^{k,2}(U)$ . The norm in  $W^{k,p}(U)$  is defined as

$$\|u\|_{W^{k,p}(U)} = \left( \sum_{|\alpha| \leq k} \int_U |D^\alpha u|^p dx \right)^{1/p}, \quad 1 \leq p < \infty,$$

and

$$\|u\|_{W^{k,\infty}(U)} = \sum_{|\alpha| \leq k} \sup_U |D^\alpha u|, \quad p = \infty.$$

## Sobolev-type inequalities

Sobolev spaces are defined in terms of weak derivatives, which brings about a natural question of how "nice in the usual sense" are functions in the Sobolev space  $W^{k,p}(U)$  for some fixed  $k$  and  $p$ . Are they continuous? Differentiable in there usual sense? This is very useful to know since norms in Sobolev spaces are in terms of integrals and are usually much easier to establish for solutions of PDEs than point-wise estimates that are required to prove continuity of differentiability of solutions.

Let us look at the following example: take  $U = \{|x| \leq 1\}$  be the unit disk in two dimensions, and set  $f(x) = x/\sqrt{x^2 + y^2}$ . This function is discontinuous, and has weak derivatives

$$\frac{\partial f(x, y)}{\partial x} = \frac{y^2}{(x^2 + y^2)^{3/2}}, \quad \frac{\partial f(x, y)}{\partial y} = -\frac{xy}{(x^2 + y^2)^{3/2}}.$$

Note that

$$\int_U \left| \frac{\partial f}{\partial x} \right| dx dy \leq \int_0^1 \frac{1}{r} r dr = 1,$$

and similarly for  $\partial f/\partial y$ . Therefore,  $f$  lies in the Sobolev space  $W^{1,1}(U)$ . On the other hand, we have

$$\int_U \left| \frac{\partial f}{\partial y} \right|^2 dx dy = \int_0^1 \int_0^{2\pi} \frac{r^4 \cos^2 \phi \sin^2 \phi}{r^6} r d\phi dr = +\infty,$$

hence  $f$  does not lie in the Sobolev space  $H^1(U)$ . This indicates that "maybe" Sobolev functions in  $H^1(U)$  are "better" than those in  $W^{1,1}(U)$ . Sobolev inequalities provide a way to quantify that.

### Gagliardo-Nirenberg inequality

Let us ask the following question; can we bound  $\|u\|_{L^q(\mathbb{R}^n)}$  in terms of  $\|\nabla u\|_{L^p(\mathbb{R}^n)}$  with some  $p$  and  $q$ ? The answer is obviously not since  $u \equiv 1$  has an infinite  $L^q$  norm for any  $1 \leq q < \infty$  but  $\nabla u \equiv 0$ . Let us restrict the question to functions  $u \in C_c^\infty(\mathbb{R}^n)$  and ask whether it is true that

$$\|u\|_{L^q(\mathbb{R}^n)} \leq C \|\nabla u\|_{L^p(\mathbb{R}^n)}, \quad \text{for all } u \in C_c^\infty(\mathbb{R}^n). \quad (4.2)$$

The constant  $C$  should not depend on the function  $u$  – hence, in particular, it would not depend on the support of  $u$ . This is very suspicious if we think of some sequence of functions  $u_m \in C_c^\infty(\mathbb{R}^n)$  that approximates the function  $u \equiv 1$  as  $m \rightarrow +\infty$ . In order to see if we have a chance, consider a family of functions  $u_\lambda(x) = u(\lambda x)$ , with  $\lambda > 0$ , and see how (4.2) holds up as we vary the constant  $\lambda$ :

$$\|u_\lambda\|_{L^q(\mathbb{R}^n)} = \left( \int_{\mathbb{R}^n} |u(\lambda x)|^q dx \right)^{1/q} = \lambda^{-n/q} \|u\|_{L^q(\mathbb{R}^n)},$$

and

$$\|\nabla u_\lambda\|_{L^p(\mathbb{R}^n)} = \left( \int_{\mathbb{R}^n} \lambda^p |\nabla u(\lambda x)|^p dx \right)^{1/p} = \lambda^{1-n/p} \|\nabla u\|_{L^p(\mathbb{R}^n)}.$$

If (4.2) holds for all  $\lambda > 0$  we should have then

$$\lambda^{-n/q} \|u\|_{L^q(\mathbb{R}^n)} \leq C \lambda^{1-n/p} \|\nabla u\|_{L^p(\mathbb{R}^n)}, \quad (4.3)$$

for all  $\lambda > 0$  and all  $u \in C_c^\infty(\mathbb{R}^n)$ . If we fix  $u$  in (4.3), we should have

$$\lambda^{n/p-n/q-1} \leq \frac{C \|\nabla u\|_{L^p(\mathbb{R}^n)}}{\|u\|_{L^q(\mathbb{R}^n)}}.$$

This is only possible if

$$\frac{n}{p} - \frac{n}{q} = 1. \quad (4.4)$$

This tells us two things: first, given  $p$ , we can only hope to prove (4.2) for

$$q = \frac{np}{n-p}, \quad (4.5)$$

and that the range of  $p$  should be restricted to  $1 \leq p < n$ .

**Theorem 4.3** (*Gagliardo-Nirenberg inequality*) *Assume that  $1 \leq p < n$ , and let  $q = np/(n-p)$ . There exists a constant  $C > 0$  so that*

$$\|u\|_{L^q(\mathbb{R}^n)} \leq C \|\nabla u\|_{L^p(\mathbb{R}^n)}, \text{ for all } u \in C_c^\infty(\mathbb{R}^n). \quad (4.6)$$

We will not prove this theorem here.

For bounded domains we have the following version.

**Theorem 4.4** *Let  $U \subset \mathbb{R}^n$  be a bounded domain with a  $C^1$ -boundary  $\partial U$ . Assume that  $1 \leq p < n$ ,  $q = np/(n-p)$ , and  $u \in W^{1,p}(U)$ . Then  $u \in L^q(U)$  and*

$$\|u\|_{L^q(U)} \leq C \|u\|_{W^{1,p}(U)}. \quad (4.7)$$

*The constant  $C$  depends only on  $p$ ,  $n$  and  $U$ .*

### Morrey's inequality

In order to understand how being in some Sobolev space and the continuity of a function are related, recall the Hölder norm a function:

$$\|u\|_{C^{0,\alpha}(U)} = \sup_{x \in U} |u(x)| + \sup_{x,y \in U} \frac{|u(x) - u(y)|}{|x - y|^\alpha}.$$

Let us see when the inequality

$$\|u\|_{C^{0,\alpha}(\mathbb{R}^n)} \leq C \|u\|_{W^{1,p}(\mathbb{R}^n)} \quad (4.8)$$

can hold. Once again, we fix  $u \in C_c^\infty(\mathbb{R}^n)$ , and consider the rescaled functions  $u_\lambda(x) = u(\lambda x)$ . The  $W^{1,p}$ -norm of the rescaled function is

$$\|u_\lambda\|_{W^{1,p}(\mathbb{R}^n)}^p = \|u_\lambda\|_{L^p}^p + \|\nabla u_\lambda\|_{L^p}^p = \lambda^{-n} \|u\|_{L^p}^p + \lambda^{p-n} \|\nabla u\|_{L^p}^p.$$

How does the Hölder norm scale with  $\lambda$ ? We have

$$\|u_\lambda\|_{C^{0,\alpha}(U)} = \sup_x |u(\lambda x)| + \sup_{x,y} \frac{|u(\lambda x) - u(\lambda y)|}{|x - y|^\alpha} = \|u\|_{C(\mathbb{R}^n)} + \lambda^\alpha F_\alpha(u).$$

Here we have set

$$F_\alpha(u) = \sup_{x,y} \frac{|u(x) - u(y)|}{|x - y|^\alpha}.$$

Therefore, for (4.8) to hold we should have, at least, that

$$\|u\|_{C(\mathbb{R}^n)} + \lambda^\alpha F_\alpha(u) \leq C(\lambda^{-n/p}\|u\|_{L^p} + \lambda^{1-n/p}\|\nabla u\|_{L^p}). \quad (4.9)$$

If  $p \leq n$  then fixing  $u$  and letting  $\lambda \rightarrow +\infty$  in (4.9) we would reach a contradiction. Hence, we have to take  $n < p \leq +\infty$ , and then we can take  $\alpha = 1 - n/p$ . This scaling analysis is confirmed by the following theorem.

**Theorem 4.5** (*Morrey's inequality*) *Assume that  $n < p \leq +\infty$ . There exists a constant  $C > 0$  that depends only on  $p$  and  $n$  so that*

$$\|u\|_{C^{0,\alpha}(\mathbb{R}^n)} \leq C\|u\|_{W^{1,p}(\mathbb{R}^n)} \quad (4.10)$$

holds for all  $u \in C^1(\mathbb{R}^n)$ , with  $\alpha = 1 - p/n$ .

### Rellich-Kondrashov compactness theorem

Gagliardo-Nirenberg inequality in a bounded domain implies that if  $1 \leq p < n$  and  $q = np/(n - p)$  then  $W^{1,p}(U)$  is a subset of  $L^q(\mathbb{R}^n)$ . Rellich-Kondrashov theorem shows that it is actually a compact subset of  $L^q(\mathbb{R}^n)$ , which is extremely important for the PDE theory.

**Definition 4.6** *Let  $X$  and  $Y$  be Banach spaces. Then  $X$  is compactly embedded in  $Y$ , written as  $X \subset\subset Y$  if*

- (i) *there exists a constant  $C > 0$  so that  $\|x\|_X \leq C\|x\|_Y$  for all  $x \in X$ .*
- (ii) *any bounded sequence  $x_n$  in  $X$  has a subsequence  $x_{n_k}$  that converges in  $Y$ .*

It is crucial, of course, that it is only required in (ii) that the subsequence converges in the space  $Y$  and not in  $X$ !

**Example.** Let  $Y = l^2$  and  $X$  be the set of all sequences  $x_n$  such that

$$\sum_{n=1}^{\infty} n^2|x_n|^2 < +\infty,$$

equipped with the norm

$$\|x\|_1 = \left( \sum_{n=1}^{\infty} n^2|x_n|^2 \right)^{1/2}.$$

A good exercise is to check that  $X$  is compactly embedded in  $Y$ . The reason is, roughly, that  $X$  "behaves as a finite-dimensional subspace of  $Y$ ". This is because if  $\|x\|_1 < 1$  then the entries  $x_n$  decay as  $|x_n| \leq 1/n^2$ , meaning that "only the few first entries of  $x$  play a role".

**Theorem 4.7** (*Rellich-Kondrashov*) *Let  $U$  be a bounded open domain in  $\mathbb{R}^n$  with a  $C^1$ -boundary  $\partial U$ . Suppose that  $1 \leq p < n$ , then  $W^{1,p}(U)$  is compactly embedded into  $L^q(U)$  for each  $1 \leq q < np/(n - p)$ .*

This theorem is crucial for construction of solutions of PDEs.

## Traces and $W_0^{1,p}(U)$ spaces

As we will be talking about PDEs with boundary conditions, we need to be able to say what it means that a function in  $W^{1,p}(U)$  (about which we do not know that it is continuous) vanishes on the boundary  $\partial U$ . This is done with the help of what is known as the trace operator. The basic estimate that makes it work is the following lemma that shows that for smooth functions their restriction to the boundary is bounded by the  $W^{1,p}$ -norm inside the domain.

**Lemma 4.8** *Let  $U$  be a bounded domain with a  $C^1$ -boundary  $\partial U$ . There exists a constant  $C > 0$  so that for all  $u \in C^1(\bar{U})$  we have*

$$\|u\|_{L^p(\partial U)} \leq C \|u\|_{W^{1,p}(U)}. \quad (4.11)$$

We will not prove this lemma but in order to understand why that can be true, consider the situation when  $U$  is a two-dimensional domain with a smooth boundary that contains the interval  $[-1, 1]$  on the  $x$ -axis. Let  $\zeta(x_1, x_2)$  be a  $C_c^\infty$  function such that  $\zeta(x, 0) = 1$  for  $x \in [-1/2, 1/2]$  and  $\zeta$  is supported inside  $U$ . Assume also, for simplicity of notation that  $u(x_1, x_2) \geq 0$  and let  $\Gamma = \{(x, 0) : -1 \leq x \leq 1\}$ . Then we have

$$\begin{aligned} \|u\|_{L^p(\Gamma)}^p &= \int_{\Gamma} |u(x, 0)|^p dx = \int_{\Gamma} \zeta(x, 0) u^p(x, 0) dx = \int_U \frac{\partial}{\partial x_2} (\zeta(x_1, x_2) u^p(x_1, x_2)) dx_1 dx_2 \\ &= \int_U \frac{\partial \zeta}{\partial x_2} u^p dx_1 dx_2 + p \int_U \zeta u^{p-1} \frac{\partial u}{\partial x_2} dx_1 dx_2 = I + II. \end{aligned}$$

For the first term we have simply

$$I \leq C \int_U |u|^p dx_1 dx_2.$$

The second can be estimated using the inequality

$$|a^{p-1}b| \leq |a|^p + |b|^p,$$

giving

$$II \leq C \int_U |u|^p dx_1 dx_2 + C \int_U |\nabla u|^p dx_1 dx_2.$$

Together, we have

$$\|u\|_{L^p(\Gamma)}^p \leq C \|u\|_{L^p(U)}^p + C \|\nabla u\|_{L^p(U)}^p = C \|u\|_{W^{1,p}}^p,$$

which is almost (4.11) (almost because  $\Gamma$  is only part of the boundary of  $U$ ). The general proof of Lemma 4.8 proceeds very similarly by making a change of variable to straighten the boundary.

Lemma 4.8 means that we can define the restriction (or trace) operator  $T : C^1(\bar{U}) \rightarrow L^p(\partial U)$  that is bounded as in (4.11). Since the set  $C^1(\bar{U})$  is dense in  $W^{1,p}(U)$  we can extend  $T$  to all functions in  $W^{1,p}(U)$  by continuity, preserving the bound

$$\|Tu\|_{L^p(\partial U)} \leq C \|u\|_{W^{1,p}(U)}. \quad (4.12)$$

Now, we can say what a zero boundary condition means for a function  $u \in W^{1,p}(U)$ .

**Definition 4.9** A function  $u \in W_0^{1,p}(U)$  if  $Tu = 0$ .

We have the following version of Gagliardo-Nirenberg inequality for bounded domains and functions in  $W_0^{1,p}(U)$ .

**Theorem 4.10** (Poincaré's inequality) Let  $U \subset \mathbb{R}^n$  be a bounded domain with a  $C^1$ -boundary  $\partial U$ . Assume that  $1 \leq p < n$ , and  $u \in W_0^{1,p}(U)$ . Then  $u \in L^q(U)$  and

$$\|u\|_{L^q(U)} \leq C \|\nabla u\|_{L^p(U)}, \quad (4.13)$$

for all  $1 \leq q \leq np/(n-p)$ . The constant  $C$  depends only on  $p, q, n$  and  $U$ .

Another useful result, also known as the Poincaré inequality is as follows.

**Theorem 4.11** (Poincaré's inequality) Let  $U \subset \mathbb{R}^n$  be a bounded domain with a  $C^1$ -boundary  $\partial U$ . There exists a constant  $C > 0$  so that for all functions  $u \in H_0^1(U)$  we have

$$\|u\|_{L^2(U)} \leq C \|\nabla u\|_{L^2(U)}. \quad (4.14)$$

Let us prove Theorem 4.11. Assume that this is false, then there exists a sequence of functions  $u_k \in H_0^1(U)$  such that

$$\|u_k\|_{L^2(U)} \geq k \|\nabla u_k\|_{L^2(U)}.$$

Consider the renormalized functions

$$v_k = \frac{u_k}{\|u_k\|_{L^2(U)}},$$

then we have

$$\|v_k\|_{L^2(U)} = 1, \quad \|\nabla v_k\|_{L^2(U)} \leq \frac{1}{k}. \quad (4.15)$$

The functions  $v_k$  are uniformly bounded in the Hilbert space  $H_0^1(U)$ , hence there exists a subsequence  $n_k$  such that  $v_{n_k}$  converges weakly to a limit  $v \in H_0^1(U)$ . Moreover, Rellich-Kondrashov theorem implies that  $v_{n_k}$  converges strongly to  $v$  in  $L^2(U)$ . It follows that  $\|v\|_{L^2(U)} = 1$  and  $\|\nabla v\|_{L^2(U)} = 0$ . The last condition implies that  $v = \text{const}$  a.e. in  $U$ . As  $v$  vanishes on  $\partial U$ , it is an easy exercise to see from Lemma 4.8 that then  $v \equiv 0$  a.e., contradicting the fact that  $\|v\|_{L^2(U)} = 1$ .

## 4.2 Weak solutions of boundary value problems

### The weak formulation for the Poisson equation

We now turn to purely PDE questions of existence and regularity of solutions to elliptic boundary value problems. Let us first consider an example of the Poisson equation

$$\begin{aligned} -\Delta u &= f \text{ in } U, \\ u &= 0 \text{ on } \partial U, \end{aligned} \quad (4.16)$$

in a bounded domain  $U$ , with a given function  $f \in L^2(U)$ . Multiplying (4.16) by a test function  $\phi \in C_c^\infty(U)$  and integrating by parts gives

$$\int_U \nabla u \cdot \nabla \phi dx = \int_U f \phi dx. \quad (4.17)$$

The left side makes sense as long as  $u$  and  $\phi$  are in  $H_0^1(U)$ , which motivates the following definition.

**Definition 4.12** We say that  $u \in H_0^1(U)$  is a weak solution of (4.16) if for any function  $v \in H_0^1(U)$  we have

$$\int_U \nabla u \cdot \nabla v dx = \int_U f v dx. \quad (4.18)$$

In order to show that (4.16) has a weak solution, let us consider  $H_0^1(U)$  equipped with the inner product

$$\langle u, v \rangle_1 = \int_U \nabla u \cdot \nabla v dx,$$

and the norm

$$\|u\|_1 = \left( \int_U |\nabla u|^2 dx \right)^{1/2}.$$

One can verify that  $H_0^1(U)$  is still a Hilbert space under this new inner product. Note, in particular, that  $\|u\|_1 = 0$  implies that  $u \equiv 0$  in  $U$  since  $u = 0$  on the boundary  $\partial U$ . The linear functional

$$A(v) = \int_U f v dx$$

is bounded on  $H_0^1(U)$  with that norm since

$$|A(v)| \leq \|f\|_{L^2} \|v\|_{L^2} \leq C \|f\|_{L^2} \|\nabla v\|_{L^2} = C \|f\|_{L^2} \|v\|_1,$$

by the Poincaré inequality (4.14). It follows now from the Riesz representation theorem for bounded functionals on Hilbert spaces that there exists an element  $w \in H_0^1(U)$  such that for any function  $v \in H_0^1(U)$  we have

$$A(v) = \langle w, v \rangle.$$

Writing this more explicitly, gives

$$\int_U f v = \int_U \nabla w \cdot \nabla v dx, \text{ for any } v \in H_0^1(U),$$

which means exactly that  $w$  is a weak solution of the Poisson equation (4.16).

### The weak formulation for general elliptic problems

We now construct weak solutions for general elliptic problems of the form

$$\begin{aligned} \mathcal{L}u &= f \text{ in } U, \\ u &= 0 \text{ on } \partial U, \end{aligned} \quad (4.19)$$

posed in a bounded domain  $U \subset \mathbb{R}^n$ . Here  $\mathcal{L}$  is an operator of the form

$$\mathcal{L}u = - \sum_{i,j=1}^n \frac{\partial}{\partial x_j} \left( a_{ij}(x) \frac{\partial u}{\partial x_i} \right) + \sum_{i=1}^m b_i(x) \frac{\partial u}{\partial x_i} + c(x)u, \quad (4.20)$$

or, in a more compact notation,

$$\mathcal{L}u = -\nabla \cdot (a(x)\nabla u) + b(x) \cdot \nabla u + c(x)u.$$

The operator  $\mathcal{L}$  in (4.20) is in its so-called divergence form. Sometimes we will also consider operators in the non-divergence form:

$$\mathcal{L}u = -\sum_{i,j=1}^n a_{ij}(x) \frac{\partial^2 u}{\partial x_i \partial x_j} + \sum_{i=1}^m b_i(x) \frac{\partial u}{\partial x_i} + c(x)u. \quad (4.21)$$

Both forms have their advantages: the divergence form is well suited for energy methods (integration by parts), while the non-divergence form is convenient for applications of the maximum principle. One can always pass from one form to another by modifying the drift  $b(x)$ . We will always assume that  $\mathcal{L}$  is elliptic, that is, there exists a constant  $c_0 > 0$  so that

$$\sum_{i,j=1}^n a_{ij}(x) \xi_i \xi_j \geq c_0 |\xi|^2, \quad (4.22)$$

for all  $\xi \in \mathbb{R}^n$ , and that the coefficients are bounded:

$$|a_{ij}(x)|, |b_i(x)|, |c(x)| \leq K, \quad (4.23)$$

with some  $K > 0$ .

We will define the notion of a weak solution for functions in  $H_0^1(U)$ , hence we need to reformulate the problem so that it would make sense for  $H_0^1(U)$  functions. In order to do that, let us take (4.19) in the divergence form (4.20), multiply it by a smooth test function  $v \in C_c^\infty$  and integrate over  $U$ :

$$\int_U \sum_{i,j=1}^n \left( a_{ij}(x) \frac{\partial u}{\partial x_i} \frac{\partial v}{\partial x_j} + v(x)(b(x) \cdot \nabla u) + c(x)uv \right) dx = \int_U f v dx. \quad (4.24)$$

The left side of (4.24) makes sense if  $u \in H_0^1(U)$ , and, in addition, the condition  $u \in H_0^1(U)$  automatically enforces the boundary condition  $u = 0$  on  $\partial U$ . In the spirit of (4.24) let us define a bilinear form

$$B(u, v) = \int_U \sum_{i,j=1}^n \left( a_{ij}(x) \frac{\partial u}{\partial x_i} \frac{\partial v}{\partial x_j} + v(x)(b(x) \cdot \nabla u) + c(x)uv \right) dx, \quad (4.25)$$

define for functions  $u, v \in H_0^1(U)$ .

**Definition 4.13** *We say that a function  $u \in H_0^1(U)$  solves the boundary value problem*

$$\begin{aligned} -\sum_{i,j=1}^n a_{ij}(x) \frac{\partial^2 u}{\partial x_i \partial x_j} + \sum_{i=1}^m b_i(x) \frac{\partial u}{\partial x_i} + c(x)u &= f \text{ in } U, \\ u &= 0 \text{ on } \partial U, \end{aligned} \quad (4.26)$$

if for any  $v \in H_0^1(U)$  we have

$$B(u, v) = \int_U f v dx. \quad (4.27)$$

Here  $B(u, v)$  is the bilinear form defined by (4.26).

### The Lax-Milgram lemma

The weak formulation (4.27) is very well suited for an application of the Lax-Milgram lemma.

**Theorem 4.14** (*Lax-Milgram lemma*) *Let  $H$  be a Hilbert space, and  $B$  be a bi-linear form on  $H$ . Assume that there exist two constant  $c_{1,2} > 0$  so that*

$$|B((u, v)| \leq c_1 \|u\|_H \|v\|_H,$$

and

$$c_2 \|u\|_H^2 \leq B(u, u).$$

Then for any bounded linear functional  $A$  on  $H$  there exists a unique  $w \in H$  such that

$$B(w, v) = A(v), \quad \text{for all } v \in H.$$

### Existence of the weak solutions

Let us now verify that the assumptions of the Lax-Milgram lemma hold for the bilinear form  $B(u, v)$  defined in (4.25).

**Lemma 4.15** *There exists  $c > 0$  so that  $|B(u, v)| \leq c \|u\|_{H_0^1(U)} \|v\|_{H_0^1(U)}$ , for all  $u, v \in H_0^1(U)$ .*

**Proof.** We simply check

$$\begin{aligned} |B(u, v)| &\leq \sum_{i,j=1}^n \|a_{ij}\|_{L^\infty(U)} \int_U |\nabla u| |\nabla v| dx + \sum_{i=1}^n \|b_i\|_{L^\infty(U)} \int_U |\nabla u| |v| dx \\ &\quad + \|c\|_{L^\infty(U)} \int_U |u| |v| dx \leq C \|u\|_{H_0^1(U)} \|v\|_{H_0^1(U)}, \end{aligned}$$

with some appropriate constant  $C$ .  $\square$

**Lemma 4.16** *There exist two constants  $c_{1,2} > 0$  so that*

$$c_1 \|u\|_{H_0^1(U)}^2 \leq B(u, u) + c_2 \|u\|_{L^2(U)}^2. \quad (4.28)$$

**Proof.** Note that, because of the ellipticity condition, we have

$$\sum_{i,j=1}^n a_{ij}(x) \frac{\partial u}{\partial x_i} \frac{\partial u}{\partial x_j} \geq c_0 |\nabla u|^2,$$

hence

$$\begin{aligned} B(u, u) &\geq c_0 \int_U |\nabla u|^2 dx - \sum_{i=1}^n \|b_i\|_{L^\infty(U)} \int_U |\nabla u| |u| dx - \|c\|_{L^\infty(U)} \int_U |u|^2 dx \\ &\geq c_0 \int_U |\nabla u|^2 dx - nK \int_U |\nabla u| |u| dx - K \int_U |u|^2 dx. \end{aligned}$$

Using the inequality

$$ab \leq \varepsilon a^2 + \frac{1}{\varepsilon} b^2,$$

gives

$$B(u, u) \geq c_0 \int_U |\nabla u|^2 dx - nK\varepsilon \int_U |\nabla u|^2 dx - \frac{nK}{\varepsilon} \int_U |u|^2 dx - K \int_U |u|^2 dx.$$

Choosing  $\varepsilon = c_0/2nK$  gives

$$B(u, u) + c_2 \int_U |u|^2 dx \geq \frac{c_0}{2} \int_U |\nabla u|^2 dx,$$

with  $c_2 = K + nK/\varepsilon$ .  $\square$

**Theorem 4.17** *There exists a number  $\gamma$  so that for each  $\mu \geq \gamma$  and each  $f \in L^2(U)$  there exists a weak solution  $u \in H_0^1(U)$  of the boundary value problem*

$$-\sum_{i,j=1}^n \frac{\partial}{\partial x_j} \left( a_{ij}(x) \frac{\partial u}{\partial x_i} \right) + \sum_{i=1}^m b_i(x) \frac{\partial u}{\partial x_i} + c(x)u + \mu u = f(x) \quad x \in U, \quad (4.29)$$

$u = 0$  on  $\partial U$ .

*In addition, this solution satisfies*

$$\|u\|_{H_0^1(U)} \leq C \|f\|_{L^2(U)}, \quad (4.30)$$

*with a constant  $C > 0$  that does not depend on the function  $f \in L^2(U)$ .*

**Proof.** Take  $c_2$  from Lemma 4.16 and consider  $\mu \geq c_2$ . Define the bilinear form

$$B_\mu(u, v) = B(u, v) + \mu(u, v).$$

Recall that

$$(u, v) = \int_U uv dx$$

is the  $L^2$  inner product. Then  $B_\mu(u, v)$  satisfies the assumptions of the Lax-Milgram lemma. Given  $f \in L^2(U)$  define also the linear functional

$$A(v) = (f, v) = \int_U f v dx.$$

The Lax-Milgram lemma now implies that there exists a unique  $u \in h_0^1(U)$  such that

$$B_\mu(u, v) = (f, v), \quad \text{for all } v \in H_0^1(U). \quad (4.31)$$

This means exactly that  $u$  is a weak solution of (4.29). In order to get the bound for  $u$  in (4.30) note that (4.31) taken with  $v = u$  gives

$$\int_U f u dx = B_\mu(u, u) = B(u, u) + \mu \|u\|_{L^2(U)}^2.$$

Using Lemma 4.16 gives

$$\int_U f u dx \geq c_1 \|u\|_{H_0^1(U)}^2 - c_2 \|u\|_{L^2(U)}^2 + \mu \|u\|_{L^2(U)}^2 \geq c_1 \|u\|_{H_0^1(U)}^2.$$

The Cauchy-Schwartz inequality implies then

$$c_1 \|u\|_{H_0^1(U)}^2 \leq \|f\|_{L^2(U)} \|u\|_{L^2(U)}.$$

Finally, the Poincaré inequality  $\|u\|_{L^2(U)} \leq C \|u\|_{H_0^1(U)}$  implies that (4.30) holds.  $\square$

Note that if the first order term in the operator  $\mathcal{L}$  vanishes:  $b_j \equiv 0$ , and the zero order term is non-negative  $c(x) \geq 0$ , then we can take  $c_2 = 0$  in Lemma 4.16, hence Theorem 4.17 applies, in this situation, to all  $\mu \geq 0$ .

### Compactness of the inverse

Let us now, once again, take  $c_2$  as in Lemma 4.16 and define the operator

$$\tilde{\mathcal{L}}u = \mathcal{L}u + c_2 u.$$

Theorem 4.17 says that the inverse operator  $\mathcal{K} = \tilde{\mathcal{L}}^{-1}$  is defined and acts on  $L^2(U)$ . Let us now re-write the equation

$$\mathcal{L}u = f, \quad u \in H_0^1(U), \tag{4.32}$$

as follows. First, (4.32) is equivalent to

$$\tilde{\mathcal{L}}u = c_2 u + f, \tag{4.33}$$

which, in turn, can be re-formulated as

$$u = \mathcal{K}(c_2 u + f), \tag{4.34}$$

or

$$\mathcal{K}u - \frac{1}{c_2} u = h, \tag{4.35}$$

where  $h = (1/c_2)\mathcal{K}f$ . The key observation in understanding when (4.35) (and hence (4.32)) has a solution, is a the following lemma.

**Lemma 4.18** *The operator  $\mathcal{K} : L^2(U) \rightarrow L^2(U)$  is a bounded linear compact operator.*

**Proof.** The fact that  $\mathcal{K}$  is a bounded operator follows from estimate (4.30). Compactness of  $\mathcal{K}$  also follows from that estimate if we recall the Rellich-Kondrashov theorem – (4.30) says that if  $\|f\|_{L^2(U)} \leq C_0$ , then  $\|\mathcal{K}f\|_{H_0^1(U)} \leq CC_0$ , and Rellich-Kondrashiov theorem says that the set  $\{\|u\|_{H_0^1(U)} \leq CC_0\}$  is a compact subset of  $L^2(U)$ . this means that  $\mathcal{K}$  maps a bounded subset of  $L^2(U)$  into a compact set and is therefore a compact operator.  $\square$

As a consequence, we know that there exist an at most countable set  $\Sigma = \{\lambda_1, \dots, \lambda_n, \dots\}$  of complex eigenvalues of the operator  $\mathcal{K}$ , and the only possible accumulation point of the eigenvalues is  $\lambda = 0$ . Moreover, for any  $\lambda \notin \Sigma$ , the equation

$$\mathcal{K}f - \lambda f = u,$$

has a unique solution for all  $\lambda \in \mathbb{C}$  and all  $u \in L^2(U)$ . Now, we have two possibilities (the Fredholm alternative): if  $1/c_2 \in \Sigma$  then (4.35) may have no solution but there exists a non-zero solution of the homogenous problem

$$\mathcal{K}u = \frac{1}{c_2}u, \quad u \neq 0.$$

This is equivalent to the fact that the problem

$$\mathcal{L}u = 0, \quad u \in H_0^1(U), \tag{4.36}$$

has non-trivial solution. The second possibility is that  $1/c_2 \notin \Sigma$ , and then (4.35) has a unique solution. An alternative way to formulate this result is as follows.

**Theorem 4.19** *There exists an at most countable set  $\Sigma \in \mathbb{R}$  so that the boundary value problem*

$$\mathcal{L}u - \lambda u = f, \tag{4.37}$$

*with the boundary condition  $u = 0$  on  $\partial U$  has a unique weak solution  $u \in H_0^1(U)$  for all  $f \in L^2(U)$ . If the set  $\Sigma$  is infinite then  $\lambda_k \rightarrow +\infty$  as  $k \rightarrow +\infty$ .*

The only claim we still need to verify here is that if  $\mathcal{L}$  has infinitely many real eigenvalues  $\lambda_k$  then  $\lambda_k \rightarrow +\infty$  as  $k \rightarrow +\infty$ . First, we have  $|\lambda_k| \rightarrow +\infty$  since  $(\mathcal{L} + c_2 I)^{-1}$  is compact. Moreover, if we have

$$\mathcal{L}u = \lambda u, \quad u \in H_0^1(U),$$

we can normalize  $u$  so that  $\|u\|_{L^2(U)} = 1$ . Then, we have

$$\begin{aligned} \lambda &= \lambda \|u\|_{L^2(U)}^2 = \int_U (\mathcal{L}u)u dx = \int_U \sum_{ij} a_{ij}(x) \frac{\partial u}{\partial x_i} \frac{\partial u}{\partial x_j} dx + \int_U \sum_{j=1}^n b_j(x) \frac{\partial u}{\partial x_j} u dx + \int_U c(x) u^2 dx \\ &\geq c_0 \int_U |\nabla u|^2 dx - C \int_U u^2(x) dx - \frac{c_0}{2} \int_U |\nabla u|^2 dx - C \int_U u^2(x) dx. \end{aligned}$$

We used above the estimate

$$\left| \int_U \sum_{j=1}^n b_j(x) \frac{\partial u}{\partial x_j} u dx \right| \leq C \int_U |\nabla u| |u| dx \leq \frac{C \cdot C}{2c_0} \int_U |u|^2 dx + \frac{C}{2C} c_0 \int_U |\nabla u|^2 dx$$

We conclude that

$$\frac{c_0}{2} \int_U |\nabla u|^2 dx \leq \lambda + K,$$

with some constant  $K > 0$  that depends only on the domain  $U$ . It follows that  $\lambda > -K$ , and we are done.  $\square$

## Regularity of solutions of elliptic equations

### The Poisson equation

Let us consider solutions of the Poisson equation

$$-\Delta u = f, \tag{4.38}$$

in the whole space  $\mathbb{R}^n$ . This equation says that "the Laplacian of  $u$  is as good as  $f$ ". Laplacian is just a linear combination of some second derivatives but it turns out that (4.38) implies, actually, that all second derivatives of  $u$  are "as good as  $f$ ". In order to see that, at least on the formal level, let us assume that  $u$  vanishes sufficiently fast at infinity, and multiply (4.38) by  $\Delta u$ :

$$\begin{aligned} \int_{\mathbb{R}^n} f^2 dx &= \int_{\mathbb{R}^n} (\Delta u)^2 dx = \sum_{i,j=1}^n \int_{\mathbb{R}^n} \frac{\partial^2 u}{\partial x_i^2} \frac{\partial^2 u}{\partial x_j^2} dx = - \sum_{i,j=1}^n \int_{\mathbb{R}^n} \frac{\partial u}{\partial x_i} \frac{\partial^3 u}{\partial x_j^2 \partial x_i} dx \\ &= \sum_{i,j=1}^n \int_{\mathbb{R}^n} \frac{\partial^2 u}{\partial x_i \partial x_j} \frac{\partial^2 u}{\partial x_i \partial x_j} dx = \sum_{i,j=1}^n \int_{\mathbb{R}^n} \left( \frac{\partial^2 u}{\partial x_i \partial x_j} \right)^2 dx. \end{aligned}$$

Therefore, for instance, if  $f \in L^2(\mathbb{R}^n)$ , and  $u \in H^1(\mathbb{R}^n)$ , it is reasonable to expect that  $u \in H^2(\mathbb{R}^n)$ . The calculation above assumes that  $u$  is sufficiently rapidly decaying at infinity and is smooth enough to justify integration by parts, making it only formal but it nevertheless captures the spirit of the elliptic regularity theory.

### Interior regularity

We consider an elliptic operator in the divergence form

$$\mathcal{L}u = - \sum_{i,j=1}^n \frac{\partial}{\partial x_i} \left( a_{ij}(x) \frac{\partial u}{\partial x_j} \right) + \sum_{i=1}^n b_i(x) \frac{\partial u}{\partial x_i} + c(x)u,$$

with the coefficients  $a_{ij}$ ,  $b_i$  and  $c$  that are uniformly bounded, and  $a_{ij}$ , as always, assumed to be uniformly positive definite. In addition, we assume here that  $a_{ij} \in C^1(U)$ .

**Theorem 4.20** *Assume that  $u \in H^1(U)$  is a weak solution of*

$$\mathcal{L}u = f, \tag{4.39}$$

*with a function  $f \in L^2(U)$ . Then for any open subset  $V \subset\subset U$  we have the estimate*

$$\|u\|_{H^2(V)} \leq C(\|f\|_{L^2(U)} + \|u\|_{L^2(U)}). \tag{4.40}$$

*The constant  $C$  depends only on  $U$ ,  $V$  and coefficients of  $\mathcal{L}$ .*

Note that we do not impose any boundary condition on  $u$  – this result holds locally inside  $U$  so the boundary condition does not matter! On the other hand, we do not allow the case  $V = U$  –  $u$  may behave very badly near the boundary  $\partial U$  unless we impose a reasonable boundary condition.

## Maximum principles

Here, we will consider elliptic operators in non-divergence form

$$\mathcal{L}u = - \sum_{i,j=1}^n a_{ij}(x) \frac{\partial^2 u}{\partial x_i \partial x_j} + \sum_{j=1}^n b_j(x) \frac{\partial u}{\partial x_j},$$

with no zero-order term. As always, we assume that the ellipticity condition holds for  $a_{ij}$ . We also assume that all coefficients  $a_{ij}$  and  $b_j$  are continuous.

**Theorem 4.21** (*Weak maximum principle*) *Assume  $u \in C^2(U) \cap C(\bar{U})$ , then, if  $\mathcal{L}u \leq 0$  in  $U$  then  $u$  attains its maximum over  $U$  on the boundary  $\partial U$ .*

If  $\mathcal{L}u \geq 0$  then looking at  $v(x) = -u(x)$  we deduce that  $u$  attains its minimum over  $U$  on the boundary  $\partial U$ .

**Proof.** First, assume that we have a strict inequality

$$\mathcal{L}u < 0 \text{ in } U. \tag{4.41}$$

If  $u(x)$  attains its maximum at an interior point  $x_0 \in U$  then  $\nabla u(x_0) = 0$ , and the Hessian matrix

$$H_{ij}(x_0) = \frac{\partial^2 u}{\partial x_i \partial x_j}(x_0)$$

is non-negative: for any vector  $\xi \in \mathbb{R}^n$  we have

$$\sum_{i,j=1}^n H_{ij}(x_0) \xi_i \xi_j \leq 0.$$

The matrix  $a_{ij}(x_0)$  is positive-definite, hence it can be decomposed as

$$a_{ij}(x) = \sum_{k=1}^n \lambda_k \xi_i^{(k)} \xi_j^{(k)}.$$

Here  $\xi^{(k)}$  is the normalized eigenvector of the matrix  $a$  corresponding to the eigenvalue  $\lambda_k > 0$ . Then we have

$$\begin{aligned} \mathcal{L}u(x_0) &= - \sum_{i,j=1}^n a_{ij}(x_0) \frac{\partial^2 u}{\partial x_i \partial x_j}(x_0) = - \sum_{i,j=1}^n a_{ij}(x_0) H_{ij}(x_0) = - \sum_{i,j,k=1}^n \lambda_k \xi_i^{(k)} \xi_j^{(k)} H_{ij}(x_0) \\ &= - \sum_{k=1}^n \lambda_k \sum_{i,j=1}^n H_{ij}(x_0) \xi_i^{(k)} \xi_j^{(k)} \geq 0, \end{aligned}$$

which is a contradiction to (4.41).

In the general case, if we only have  $\mathcal{L}u \leq 0$ , set

$$v(x) = u(x) + \varepsilon e^{\lambda x_1},$$

then

$$\mathcal{L}v = \mathcal{L}u + \varepsilon \mathcal{L}(e^{\lambda x_1}) \leq \varepsilon \mathcal{L}(e^{\lambda x_1}) = \varepsilon[-\lambda^2 a_{11} + \lambda b_1]e^{\lambda x_1}.$$

As the matrix  $a_{ij}$  is positive-definite,  $a_{11}(x) \geq c_0 > 0$  in  $U$ . Therefore, if we choose  $\lambda > 0$  sufficiently large (independent of  $\varepsilon > 0$ ), we will have  $\mathcal{L}v < 0$  and thus  $v(x)$  would attain its maximum on the boundary  $\partial U$ :

$$v(x) \leq M_\varepsilon = \max_{y \in \partial U} v(y), \quad \text{for all } x \in U.$$

However, for  $\lambda > 0$  fixed we have

$$u(x) = v(x) - \varepsilon e^{\lambda x_1} \leq v(x) \leq M_\varepsilon,$$

for all  $\varepsilon > 0$ . Letting  $\varepsilon \rightarrow 0$  we conclude that

$$u(x) \leq M = \max_{y \in \partial U} u(y),$$

and we are done.  $\square$

Let us now consider domains satisfying the interior ball condition: for every  $x \in \partial U$  there exists a ball  $B \subset U$  such that  $x \in \partial B$  – this condition is automatically satisfied if  $\partial U$  is  $C^2$ .

**Lemma 4.22** (*Hopf's Lemma*) *Assume that  $u \in C^2(U) \cap C(\bar{U})$  and  $\mathcal{L}u \leq 0$  in  $U$ . Suppose in addition, that  $U$  satisfies the interior ball condition and there exists a point  $x_0 \in \partial U$  such that  $u(x_0) > u(x)$  for all  $x \in U$ , then*

$$\frac{\partial u}{\partial \nu}(x_0) > 0.$$

**Proof.** Let us assume that the ball that satisfies the interior ball condition at  $x_0$  is  $B(y, r)$ . We will choose  $\lambda > 0$  so that the function

$$v(x) = e^{-\lambda|x-y|^2} - e^{-\lambda r^2}$$

would satisfy

$$\mathcal{L}v \leq 0. \tag{4.42}$$

Let us compute

$$\begin{aligned} \mathcal{L}v &= e^{-\lambda|x-y|^2} \sum_{i,j=1}^n a_{ij}(-4\lambda^2(x_i - y_i)(x_j - y_j) + 2\lambda\delta_{ij}) - e^{-\lambda|x-y|^2} \sum_{j=1}^n b_j(-2\lambda(x_j - y_j)) \\ &\leq e^{-\lambda|x-y|^2} \left[ -c_0\lambda^2|x-y|^2 + 2\lambda \sum_{i=1}^n a_{ii} + 2\lambda|b||x-y| \right]. \end{aligned}$$

If we choose now  $\lambda > 0$  sufficiently large (depending on  $r$ ), we have, in the annulus  $D = \{r/2 < |x-y| < r\}$ :

$$\mathcal{L}v \leq e^{-\lambda|x-y|^2} \left[ -c_0\lambda^2 \frac{r^2}{4} + 2\lambda \sum_{i=1}^n a_{ii} + 2\lambda|b|r \right] \leq 0.$$

Next, note that since  $u(x_0) > u(x)$  for all  $x \in U$ , hence there exists  $\varepsilon > 0$  so that in the smaller ball we have

$$u(x_0) \geq u(x) + \varepsilon v(x) \text{ for } x \in B(y, r/2).$$

We also have

$$u(x_0) \geq u(x) + \varepsilon v(x) \text{ for } x \in \partial B(y, r),$$

simply because  $v = 0$  on  $\partial B(y, r)$ . Now, we are ready to apply the weak maximum principle in the annulus  $D$  to the function

$$p(x) = u(x) + \varepsilon v(x) - u(x_0).$$

This function satisfies  $\mathcal{L}p(x) \leq 0$  in  $D$ , and  $p(x) \leq 0$  on  $\partial R$ . Hence,  $p(x) \leq 0$  in  $D$  by the weak maximum principle. As  $p(x_0) = 0$ , we have

$$\frac{\partial p}{\partial \nu}(x_0) \geq 0.$$

It follows that

$$\frac{\partial u}{\partial \nu}(x_0) \geq -\varepsilon \frac{\partial v}{\partial \nu}(x_0) > 0,$$

and we are done.  $\square$

**Theorem 4.23** (*Strong maximum principle*) *Assume that  $U$  is a connected open bounded domain,  $u \in C^2(U) \cap \bar{U}$  and  $\mathcal{L}u \leq 0$  in  $U$ . If  $u$  attains its maximum at an interior point of  $U$  then  $u \equiv \text{const}$  in  $U$ .*

**Proof.** Let  $M = \max_{\bar{U}} u(x)$  and consider the set  $C = \{x \in U : u(x) = M\}$ . Assume that  $u \not\equiv M$ , then  $V = \{x \in U : u(x) < M\}$  is a non-empty open set. Choose a point  $y \in V$  such that  $\text{dist}(y, C) < \text{dist}(y, \partial U)$ . Let  $B$  be the largest ball centered at  $y$  that lies in  $V$ . Then there exists a point  $z \in C$  such that  $z \in \partial B$ . Then  $V$  satisfies the interior ball condition at  $z$ , hence  $\partial u / \partial \nu(z) > 0$ . But this is a contradiction –  $u$  attains its interior maximum at  $z$  whence  $\nabla u(z) = 0$ .  $\square$

## 5 Homogenization of elliptic equations

Homogenization theory deals with the following issue: consider a partial differential equation with inhomogeneous coefficients, say,

$$\begin{aligned} -\nabla \cdot (a(x)\nabla \phi) &= f(x), \text{ in } \Omega, \\ \phi &= 0 \text{ on } \partial\Omega. \end{aligned} \tag{5.1}$$

If  $a(x)$  is “highly non-uniform” then numerical solution of (5.1) may be very expensive. Moreover, in practice, in cases when  $a(x)$  is oscillatory, we often do not have the precise measurements of  $a(x)$  – it is usually unknown or known only approximately, hence we can not even attempt to solve (5.1) numerically – we do not know the coefficients! An engineering

approach is often to replace the highly heterogenous diffusivity  $a(x)$  by a constant matrix  $\bar{a}$ . That is, (5.1) is replaced by

$$\begin{aligned} -\nabla \cdot (\bar{a} \nabla \bar{\phi}) &= f(x), \text{ in } \Omega, \\ \bar{\phi} &= 0 \text{ on } \partial\Omega, \end{aligned} \tag{5.2}$$

which is a (much simpler!) homogeneous problem. The basic mathematical question is when this is justified, that is, for what class of  $a(x)$  we can find an “effective diffusivity”  $\bar{a}$  so that the solution  $\phi(x)$  of the original problem (5.1) is close to  $\bar{\phi}(x)$ , solution of the homogenized problem (5.2).

This issues is especially important if the basic PDE problem is not as regular as (5.1), which is a very nice elliptic problem with all the regularizing properties that come with ellipticity. If the underlying PDE is less regularizing then find the exact numerical solution with the oscillatory coefficients may be numerically nearly impossible. Moreover, the small scale oscillations in the solution may be of no interest to us – all we need to find are the “large scale” features of the solution. The latter can be well captured by an effective homogenized problem meaning that if we choose the effective parameters correctly then the relevant information about the true solution is well approximated by the solution of the effective homogeneous problem that is easy to obtain numerically.

This program usually can be carried out when the coefficients in the PDE, such as the diffusivity matrix  $a(x)$  in (5.1) have some “spatially homogeneous” structure and when the domain  $\Omega$  is much larger than the scale of variations of  $a(x)$ , so that  $\Omega$  consists of many sub-domains where  $a(x)$  behaves similarly. The simplest example is when  $a(x)$  is periodic, and  $\Omega$  contains many period cells of  $a(x)$ . A way to formalize this relation is to assume that  $a(x)$  has the form  $a_\varepsilon(x) = a_0(x/\varepsilon)$  where  $a_0(x)$  is a 1-periodic function in all directions. That is,  $a_\varepsilon(x)$  is  $\varepsilon$ -periodic in all directions. Fixing  $\Omega$  and letting  $\varepsilon \rightarrow 0$  corresponds then to the situation when  $\Omega$  contains  $N_\varepsilon = (1/\varepsilon)^n$  period cells of  $a_\varepsilon$ . The mathematical problem is then as follows: consider solutions of

$$\begin{aligned} -\nabla \cdot (a(\frac{x}{\varepsilon}) \nabla \phi_\varepsilon) &= f(x), \text{ in } \Omega, \\ \phi_\varepsilon &= 0 \text{ on } \partial\Omega, \end{aligned} \tag{5.3}$$

with a periodic function  $a(y)$ . Then we need to find a homogenized matrix  $\bar{a}$  that does not depend on  $x$  so that, for all  $f(x)$  in some class, solutions of (5.3) and the effective problem

$$\begin{aligned} -\nabla \cdot (\bar{a} \nabla \bar{\phi}) &= f(x), \text{ in } \Omega, \\ \bar{\phi} &= 0 \text{ on } \partial\Omega, \end{aligned} \tag{5.4}$$

satisfy  $\|\phi_\varepsilon - \bar{\phi}\| \rightarrow 0$  as  $\varepsilon \rightarrow 0$ , in some appropriate norm. A much more difficult problem is to address the same question when the function  $a(y)$  is a random field that is statistically homogeneous in space.

## 5.1 One-dimensional elliptic homogenization

The simplest periodic problem where periodic homogenization can be done is the one-dimensional elliptic problem, where all computations are very explicit – this is essentially the only case where everything can be done by hand.

Let us first recall a simple fact that holds in any number of dimensions.

**Theorem 5.1** Let  $a(x)$  be a continuous 1-periodic function and set  $a_\varepsilon(x) = a(x/\varepsilon)$ . Given any compact set  $\Omega$ ,  $a_\varepsilon(x)$  converge as  $\varepsilon \rightarrow 0$  weakly in  $L^2(\Omega)$  to the constant function

$$\bar{a} = \int_{\mathbb{T}^n} a(y) dy.$$

Here  $\mathbb{T}^n = [0, 1]^n$  is the period cell of  $a(y)$ .

**Proof.** Let  $\psi(x)$  be a smooth compactly supported test function, then we have

$$\int \psi(x) a\left(\frac{x}{\varepsilon}\right) dx = \sum_{k \in \mathbb{Z}^n} \hat{a}_k \int e^{2\pi i \xi \cdot x + 2\pi i k \cdot x/\varepsilon} \hat{\psi}(\xi) dx d\xi = \sum_{k \in \mathbb{Z}^n} \hat{a}_k \hat{\psi}\left(-\frac{k}{\varepsilon}\right) \rightarrow a_0 \hat{\psi}(0) = \bar{a} \int \psi(x) dx,$$

as  $\varepsilon \rightarrow 0$ . Here we have used the Lebesgue dominated convergence theorem and also defined the Fourier coefficients

$$a_k = \int_{\mathbb{T}^n} e^{-ik \cdot x} a(x) dx,$$

and the Fourier transform of  $\psi$ :

$$\hat{\psi}(\xi) = \int_{\mathbb{R}^n} e^{-i\xi \cdot x} \psi(x) dx.$$

As the functions  $a(x/\varepsilon)$  are uniformly bounded in  $L^2(\Omega)$  (because they are in  $L^\infty(\Omega)$  and  $\Omega$  is a bounded domain), weak convergence in  $L^2(\Omega)$  follows from the density of smooth compactly supported functions in  $L^2(\Omega)$ .

We now turn to the one-dimensional elliptic problem

$$-\frac{d}{dx} \left( a\left(\frac{x}{\varepsilon}\right) \frac{d\phi_\varepsilon}{dx} \right) = f(x), \quad (5.5)$$

with the boundary condition  $\phi_\varepsilon(0) = \phi_\varepsilon(1) = 0$ . There are several ways to homogenize (5.5) that are all interesting in their own right.

### Homogenization via an explicit solution

The most obvious one is to solve (5.5) explicitly. Integrating this equation we have

$$-a\left(\frac{x}{\varepsilon}\right) \frac{d\phi_\varepsilon}{dx} = B_1^\varepsilon + \int_0^x f(y) dy.$$

Dividing by  $a(x/\varepsilon)$  and integrating again, using the boundary condition at  $x = 0$  gives

$$-\phi_\varepsilon(x) = B_1^\varepsilon \int_0^x \frac{dy}{a(y/\varepsilon)} + \int_0^x \left( \int_0^y \frac{f(z)}{a(y/\varepsilon)} dz \right) dy = B_1^\varepsilon \int_0^x \frac{dy}{a(y/\varepsilon)} + \int_0^x \left( \int_z^x \frac{f(z)}{a(y/\varepsilon)} dy \right) dz.$$

The boundary condition at  $x = 1$  implies that the constant  $B_1^\varepsilon$  is determined by

$$B_1^\varepsilon \int_0^1 \frac{dy}{a(y/\varepsilon)} + \int_0^1 f(z) \left( \int_z^1 \frac{1}{a(y/\varepsilon)} dy \right) dz = 0.$$

Passing to the limit  $\varepsilon \rightarrow 0$  above, and using Theorem 5.1 we get that

$$\bar{B} = \lim_{\varepsilon \rightarrow 0} B_1^\varepsilon$$

satisfies

$$\bar{B} \frac{1}{\bar{a}_1} + \frac{1}{\bar{a}_1} \int_0^1 (1-z)f(z)dz = 0. \quad (5.6)$$

Here we have defined  $\bar{a}_1$  by

$$\frac{1}{\bar{a}_1} = \int_0^1 \frac{dy}{a(y)}. \quad (5.7)$$

Going back to  $\phi_\varepsilon(x)$ , and using the above limit for  $B_1^\varepsilon$ , as well as Theorem 5.1 once again, we conclude that  $\phi_\varepsilon(x)$  converges as  $\varepsilon \rightarrow 0$  to

$$\bar{\phi}(x) = \frac{x}{\bar{a}_1} \int_0^1 (1-z)f(z)dz - \frac{1}{\bar{a}_1} \int_0^x (x-z)f(z)dz. \quad (5.8)$$

it is easy to verify that the function  $\bar{\phi}(x)$  satisfies

$$-\frac{d}{dx} \left( \bar{a}_1 \frac{d\bar{\phi}}{dx} \right) = f(x), \quad (5.9)$$

with the boundary condition  $\bar{\phi}(0) = \bar{\phi}(1) = 0$ . Therefore, the homogenized diffusion coefficient in this case is given not by the spatial average of the function  $a(y)$ , as one might expect naively but by (5.7).

### One-dimensional homogenization Franco-Italian style

The second way to homogenize (5.5) is quite elegant. The solution of

$$-\frac{d}{dx} \left( a\left(\frac{x}{\varepsilon}\right) \frac{d\phi_\varepsilon}{dx} \right) = f(x), \quad (5.10)$$

with  $\phi_\varepsilon(0) = \phi_\varepsilon(1) = 0$  satisfies the identity

$$\int_0^1 a\left(\frac{x}{\varepsilon}\right) \left| \frac{d\phi_\varepsilon}{dx} \right|^2 dx = \int_0^1 f(x)\phi_\varepsilon(x)dx \quad (5.11)$$

obtained by multiplying (5.10) by  $\phi_\varepsilon$  and integrating by parts. The Poincaré inequality for functions that vanish at  $x = 0$  and  $x = 1$  is

$$\int_0^1 |\phi(x)|^2 dx \leq \int_0^1 |\phi'(x)|^2 dx. \quad (5.12)$$

Using this in (5.11), together with the ellipticity condition

$$0 < c_1 \leq a(y) \leq c_2 < +\infty, \text{ for all } y \in [0, 1],$$

gives

$$c_1 \int_0^1 |\phi'_\varepsilon(x)|^2 dx \leq \|f\|_2 \|\phi_\varepsilon\|_2 \leq \|f\|_2 \|\phi'_\varepsilon\|_2, \quad (5.13)$$

from which we conclude that

$$\|\phi'_\varepsilon\|_2 \leq \frac{1}{c_1} \|f\|_2. \quad (5.14)$$

As a consequence of this uniform bound in  $H_0^1[0, 1]$  on  $\phi_\varepsilon$ , we deduce that, after extraction of a subsequence  $\varepsilon_k \rightarrow 0$ , the functions  $\phi_{\varepsilon_k}$  converge weakly in  $H_0^1(0, 1]$ , and strongly in  $L^2[0, 1]$  to a function  $\bar{\phi} \in H_0^1[0, 1]$ . Moreover, by the same token the sequence

$$v_{\varepsilon_k}(x) = a\left(\frac{x}{\varepsilon_k}\right) \frac{d\phi_{\varepsilon_k}}{dx} \quad (5.15)$$

is bounded in  $L^2[0, 1]$  and satisfies

$$-\frac{dv_{\varepsilon_k}}{dx} = f,$$

whence is uniformly bounded in  $H^1[0, 1]$ . Therefore, it converges weakly (possibly after extracting a subsequence) in  $H^1[0, 1]$  and strongly in  $L^2[0, 1]$  to a limit  $\bar{v}(x)$  that satisfies

$$-\frac{d\bar{v}}{dx} = f(x). \quad (5.16)$$

In order to relate  $\bar{\phi}$  and  $\bar{v}$  consider (5.15) written as

$$\frac{d\phi_{\varepsilon_k}}{dx} = \frac{1}{a\left(\frac{x}{\varepsilon_k}\right)} v_{\varepsilon_k}(x). \quad (5.17)$$

Passing to the limit  $\varepsilon_k \rightarrow 0$  and using the strong convergence of  $v_{\varepsilon_k}$  in  $L^2[0, 1]$  paired with the weak convergence of the sequence  $1/a(x/\varepsilon_k)$  to  $1/\bar{a}_1$  we get

$$\frac{d\bar{\phi}}{dx} = \frac{1}{\bar{a}_1} \bar{v}. \quad (5.18)$$

Putting together (5.16) and (5.18) gives

$$-\frac{d}{dx} \left( \bar{a}_1 \frac{d\bar{\phi}}{dx} \right) = f(x), \quad (5.19)$$

and the boundary condition  $\bar{\phi}(0) = \bar{\phi}(1) = 0$  is encoded in the fact that  $\bar{\phi} \in H_0^1[0, 1]$ .

## Homogenization via multiple scales expansions

The last method is probably the most popular. Once again, we start with

$$-\frac{d}{dx} \left( a\left(\frac{x}{\varepsilon}\right) \frac{d\phi_\varepsilon}{dx} \right) = f(x), \quad (5.20)$$

with  $\phi_\varepsilon(0) = \phi_\varepsilon(1) = 0$ . Now, we look for the solution in the form of an asymptotic series

$$\phi_\varepsilon(x) = \bar{\phi}\left(x, \frac{x}{\varepsilon}\right) + \varepsilon \phi_1\left(x, \frac{x}{\varepsilon}\right) + \varepsilon^2 \phi_2\left(x, \frac{x}{\varepsilon}\right) + \dots \quad (5.21)$$

We assume here that the functions  $\bar{\phi}(x, y)$ ,  $\phi_1(x, y)$  and  $\phi_2(x, y)$  are periodic in the “fast” variable  $y \in [0, 1]$ . The plan is to insert this series into (5.20) and equate the coefficients at various powers of  $\varepsilon$ . We have the following rule:

$$\frac{d}{dx}u(x, \frac{x}{\varepsilon}) = \frac{\partial u}{\partial x}(x, \frac{x}{\varepsilon}) + \frac{1}{\varepsilon} \frac{\partial u}{\partial y}(x, \frac{x}{\varepsilon}),$$

hence (5.20) can be formally written as

$$-\left(\frac{d}{dx} + \frac{1}{\varepsilon} \frac{d}{dy}\right) \left(a(y) \left(\frac{d}{dx} + \frac{1}{\varepsilon} \frac{d}{dy}\right) (\bar{\phi}(x, y) + \varepsilon \phi_1(x, y) + \varepsilon^2 \phi_2(x, y) + \dots)\right) = f(x). \quad (5.22)$$

The term at  $\varepsilon^{-2}$  is

$$-\frac{d}{dy} \left(a(y) \frac{d\bar{\phi}}{dy}\right) = 0. \quad (5.23)$$

As the function  $\bar{\phi}(x, y)$  is 1-periodic in  $y$  we deduce from (5.23) that this function does not depend on  $y$ :

$$\bar{\phi} = \bar{\phi}(x). \quad (5.24)$$

The term of the order  $\varepsilon^{-1}$  in (5.22) is

$$-\frac{d}{dx} \left(a(y) \frac{d\bar{\phi}}{dy}\right) - \frac{d}{dy} \left(a(y) \frac{d\bar{\phi}}{dx}\right) - \frac{d}{dy} \left(a(y) \frac{d\phi_1}{dy}\right) = 0. \quad (5.25)$$

Taking into account (5.24) we obtain

$$-\frac{d}{dy} \left(a(y) \frac{d\phi_1}{dy}\right) = a'(y) \frac{d\bar{\phi}}{dx}. \quad (5.26)$$

Therefore, the function  $\phi_1(x, y)$  can be written as a product

$$\phi_1(x, y) = \chi(y) \frac{d\bar{\phi}}{dx}, \quad (5.27)$$

with a periodic function  $\chi(y)$  that satisfies

$$-\frac{d}{dy} \left(a(y) \frac{d\chi}{dy}\right) = a'(y). \quad (5.28)$$

The function  $\chi(y)$  is called the corrector, this equation is known as “the cell problem”. It follows that

$$-a(y) \frac{d\chi}{dy} = C_0 + a(y). \quad (5.29)$$

The constant  $C_0$  is determined by the requirement that  $\chi$  is periodic in  $y$ : integrating (5.29) in  $y$  gives, for a periodic  $\chi$ :

$$C_0 = -\left(\int_0^1 \frac{dy}{a(y)}\right)^{-1} = -\bar{a}_1. \quad (5.30)$$

The function  $\chi(y)$  is then determined up to a constant  $C_1$ :

$$\chi(y) = C_1 + \bar{a}_1 \int_0^y \frac{dz}{a(z)} - z. \quad (5.31)$$

The terms of order  $\varepsilon^0$  in (5.22) are

$$-\frac{d}{dx} \left( a(y) \frac{d\bar{\phi}}{dx} \right) - \frac{d}{dx} \left( a(y) \frac{d\phi_1}{dy} \right) - \frac{d}{dy} \left( a(y) \frac{d\phi_1}{dx} \right) - \frac{d}{dy} \left( a(y) \frac{d\phi_2}{dy} \right) = f(x). \quad (5.32)$$

Integrating in  $y$  over  $[0, 1]$  gives:

$$-\bar{a} \frac{d^2 \bar{\phi}}{dx^2} - \left( \int_0^1 a(y) \chi'(y) dy \right) \frac{d^2 \bar{\phi}}{dx^2} = f(x), \quad (5.33)$$

with

$$\bar{a} = \int_0^1 a(y) dy.$$

Note that

$$\int_0^1 a(y) \chi'(y) dy = - \int_0^1 (C_0 + a(y)) dy = \bar{a}_1 - \bar{a}.$$

Using this in (5.33) we obtain the homogenized equation

$$-\bar{a}_1 \frac{d^2 \bar{\phi}}{dx^2} = f(x). \quad (5.34)$$

Ironically, while this method seems the clumsiest of the techniques we used to derive the homogenized problem, it is the most effective approach in dimensions greater than one where explicit solutions are not readily available. Moreover, the asymptotic expansion

$$\phi_\varepsilon(x) = \bar{\phi}(x) + \varepsilon \phi_1(x, \frac{x}{\varepsilon}) + \varepsilon^2 \phi_2(x, \frac{x}{\varepsilon}) + \dots$$

also explains what happens to  $\phi'_\varepsilon(x)$ : differentiating the expansion we get

$$\phi'_\varepsilon(x) = \bar{\phi}'(x) + \frac{\partial \phi_1(x, x/\varepsilon)}{\partial y} + \varepsilon \frac{\partial \phi_1(x, x/\varepsilon)}{\partial x} + \varepsilon \frac{\partial \phi_2(x, x/\varepsilon)}{\partial y} + \dots \quad (5.35)$$

Recalling (5.27) we see that

$$\phi'_\varepsilon(x) = \bar{\phi}'(x) + \chi'\left(\frac{x}{\varepsilon}\right) \bar{\phi}'(x) + O(\varepsilon), \quad (5.36)$$

and, indeed, as we will see later in higher dimensions, we have

$$\|\phi'_\varepsilon(x) - \bar{\phi}'(x) - \chi'\left(\frac{x}{\varepsilon}\right) \bar{\phi}'(x)\|_{L^2} \rightarrow 0, \quad (5.37)$$

as  $\varepsilon \rightarrow 0$ . This means, in particular, that while  $\phi_\varepsilon$  converges to  $\bar{\phi}$  in  $L^2[0, 1]$ , it does not converge to  $\bar{\phi}$  in  $H_0^1(\Omega)$ .

## 5.2 Homogenization in dimensions $d \geq 2$

### The multiple scales expansion

We now consider elliptic homogenization in dimension  $d \geq 2$ . We consider the boundary value problem

$$\begin{aligned} -\nabla \cdot \left( a\left(\frac{x}{\varepsilon}\right) \nabla \phi_\varepsilon \right) &= f(x), \text{ in } \Omega \\ \phi &= 0 \text{ on } \partial\Omega. \end{aligned} \quad (5.38)$$

It is important to note that, as in the one-dimensional case, the smooth bounded domain  $\Omega \subset \mathbb{R}^n$  is fixed and the scale  $\varepsilon$  of the oscillations of the diffusivity matrix  $a(x/\varepsilon)$  will be sent to zero. We assume that  $a(y)$  is a smooth matrix-valued function in all 1-periodic variables  $x_j, j = 1, \dots, n$ , and denote by  $\mathbb{T}^n$  the unit  $n$ -dimensional torus:  $\mathbb{T}^n = [0, 1]^n$ . We also assume that  $a_{ij}(x)$  is uniformly elliptic: there exists a constant  $c > 0$  so that for all  $y \in \mathbb{T}^n$  and all  $\xi \in \mathbb{R}^n$  we have

$$c_0 |\xi|^2 \leq \sum_{i,j=1}^n a_{ij}(x) \xi_i \xi_j \leq c_0^{-1} |\xi|^2. \quad (5.39)$$

The “explicit” solution method we used in one dimension does not apply any more, and we first describe the approach via the multiple series expansion that is somewhat tedious but nearly universally effective. We look for the solution in the form of an asymptotic series

$$\phi_\varepsilon(x) = \bar{\phi}\left(x, \frac{x}{\varepsilon}\right) + \varepsilon \phi_1\left(x, \frac{x}{\varepsilon}\right) + \varepsilon^2 \phi_2\left(x, \frac{x}{\varepsilon}\right) + \dots \quad (5.40)$$

Here  $x \in \Omega$  is the “slow variable”, and  $y \in \mathbb{T}^n$  is the “fast variable”. As in one dimension, we assume that the functions  $\bar{\phi}(x, y)$ ,  $\phi_1(x, y)$  and  $\phi_2(x, y)$  are periodic in the “fast” variable  $y \in \mathbb{T}^n$ . We will insert this series into (5.38) and equate the coefficients at various powers of  $\varepsilon$  with the following rule:

$$\nabla_x u\left(x, \frac{x}{\varepsilon}\right) = \nabla_x u\left(x, \frac{x}{\varepsilon}\right) + \frac{1}{\varepsilon} \nabla_y u\left(x, \frac{x}{\varepsilon}\right),$$

hence (5.38) can be formally written as

$$-\left( \nabla_x + \frac{1}{\varepsilon} \nabla_y \right) \cdot \left( a(y) \left( \nabla_x + \frac{1}{\varepsilon} \nabla_y \right) \left( \bar{\phi}(x, y) + \varepsilon \phi_1(x, y) + \varepsilon^2 \phi_2(x, y) + \dots \right) \right) = f(x). \quad (5.41)$$

The term at  $\varepsilon^{-2}$  in (5.41) is

$$-\nabla_y \cdot \left( a(y) \nabla_y \bar{\phi} \right) = 0. \quad (5.42)$$

As the function  $\bar{\phi}(x, y)$  is 1-periodic in  $y$ , and  $a(y)$  is uniformly positive and bounded, we deduce from (5.42) that this function does not depend on  $y$ :

$$\bar{\phi} = \bar{\phi}(x). \quad (5.43)$$

In some sense, this is the signature of homogenization: the leading order term does not oscillate on the fast scale.

The term of the order  $\varepsilon^{-1}$  in (5.41) is

$$-\nabla_x \cdot (a(y)\nabla_y \bar{\phi}) - \nabla_y \cdot (a(y)\nabla_x \bar{\phi}) - \nabla_y \cdot (a(y)\nabla_y \phi_1) = 0. \quad (5.44)$$

Taking into account (5.43) we obtain

$$-\nabla_y \cdot (a(y)\nabla_y \phi_1) = \nabla_y a(y) \cdot \nabla_x \bar{\phi} = \sum_{j,k=1}^n \frac{\partial a_{jk}}{\partial y_j} \frac{\partial \bar{\phi}}{\partial x_k}. \quad (5.45)$$

The right side has a product structure, and the left side involves only an operator in the fast variable  $y$  – hence solution of (5.45) may be decomposed as

$$\phi_1(x, y) = \sum_{k=1}^n \chi_k(y) \frac{\partial \bar{\phi}(x)}{\partial x_k}. \quad (5.46)$$

The periodic function  $\chi(y) = (\chi_1(y), \dots, \chi_n(y))$  satisfies

$$-\nabla_y \cdot (a(y)\nabla_y \chi_k) = \sum_{j=1}^n \frac{\partial a_{jk}(y)}{\partial y_j}. \quad (5.47)$$

The function  $\chi(y)$  is called the corrector. This equation known as “the cell problem”, or the “corrector equation”, and its analogs in other homogenization problems are the key to the homogenization procedure. Unlike in one dimension, the cell problem does not have an explicit solution. Note that the function  $\chi_j(y)$  is determined only up to an additive constant  $C_j$ . Therefore, our final answer should be invariant with respect to an addition of an arbitrary constant vector to  $\chi(y)$ .

The terms of order  $\varepsilon^0$  in (5.41) are

$$-\nabla_x \cdot (a_0 \nabla_x \bar{\phi}) - \nabla_x \cdot (a(y)\nabla_y \phi_1) - \nabla_y \cdot (a(y)\nabla_x \phi_1) - \nabla_y \cdot (a(y)\nabla_y \phi_2) = f(x). \quad (5.48)$$

Integrating in  $y$  over  $\mathbb{T}^n$  gives, using expression (5.46) for  $\phi_1(x, y)$ :

$$-\nabla_x \cdot (a_0 \nabla_x \bar{\phi}) - \sum_{k,m,j=1}^n \frac{\partial}{\partial x_k} \left( \int_{\mathbb{T}^n} a_{km}(y) \frac{\partial \chi_j(y)}{\partial y_m} dy \right) \frac{\partial \bar{\phi}}{\partial x_j} = f(x), \quad (5.49)$$

with the matrix

$$a_0 = \int_{\mathbb{T}^n} a(y) dy.$$

Therefore, we have obtained the following homogenized problem

$$\begin{aligned} -\nabla \cdot (\bar{a} \nabla \bar{\phi}) &= f(x) \text{ in } \Omega, \\ \bar{\phi} &= 0 \text{ on } \partial\Omega. \end{aligned} \quad (5.50)$$

The effective diffusivity matrix is given by

$$\bar{a}_{kj} = \int_{\mathbb{T}^n} \left( a_{kj}(y) + \sum_{m=1}^n a_{km}(y) \frac{\partial \chi_j(y)}{\partial y_m} \right) dy, \quad (5.51)$$

with the functions  $\chi_j$  that satisfy the cell problem (5.47). The matrix  $\bar{a}$  depends only on  $\nabla\chi_j$  and thus does not change if we add an arbitrary constant to  $\chi_j$ .

As in one dimension, the asymptotic expansion

$$\phi_\varepsilon(x) = \bar{\phi}(x) + \varepsilon\phi_1(x, \frac{x}{\varepsilon}) + \varepsilon^2\phi_2(x, \frac{x}{\varepsilon}) + \dots$$

explains what happens to  $\nabla\phi_\varepsilon(x)$ . Differentiating the expansion we get

$$\frac{\partial\phi_\varepsilon(x)}{\partial x_j} = \frac{\partial\bar{\phi}(x)}{\partial x_j} + \frac{\partial\phi_1(x, x/\varepsilon)}{\partial y_j} + \varepsilon\frac{\partial\phi_1(x, x/\varepsilon)}{\partial x_j} + \varepsilon\frac{\partial\phi_2(x, x/\varepsilon)}{\partial y_j} + \dots \quad (5.52)$$

We see that

$$\frac{\partial\phi_\varepsilon(x)}{\partial x_j} = \frac{\partial\bar{\phi}(x)}{\partial x_j} + \sum_{m=1}^n \frac{\partial\chi_m(x/\varepsilon)}{\partial y_j} \frac{\partial\bar{\phi}(x)}{\partial x_m} + O(\varepsilon). \quad (5.53)$$

As in one dimension, this means, in particular, that while  $\phi_\varepsilon$  converges to  $\bar{\phi}$  in  $L^2[0, 1]$ , it does not converge to  $\bar{\phi}$  in  $H_0^1(\Omega)$ .

### Properties of the effective diffusion matrix

Let us now establish some basic properties of the effective diffusion matrix  $\bar{a}$ : we will show that it is symmetric and positive-definite so that the homogenized problem is well-posed. It is helpful to introduce the bilinear form

$$L(\phi, \psi) = \int_{\mathbb{T}^n} (a(y)\nabla\phi(y) \cdot \nabla\psi(y))dy, \quad (5.54)$$

defined for smooth periodic functions  $\phi, \psi \in C^1(\mathbb{T}^n)$ . The key observation is the following.

**Proposition 5.2** *The weak form of the cell problem*

$$-\nabla_y \cdot (a(y)\nabla_y\chi_k) = \sum_{j=1}^n \frac{\partial a_{jk}(y)}{\partial y_j} \quad (5.55)$$

is

$$L(\phi, \chi_k + y_k) = 0, \quad \text{for all } \phi \in C^1(\mathbb{T}^n), k = 1, \dots, n. \quad (5.56)$$

**Proof.** Fix some  $j$  and let  $e_j$  be the unit vector in the direction of  $x_j$ , then the cell problem (5.55) has the form

$$-\nabla_y \cdot (a(y)\nabla_y\chi_k) = \sum_{j=1}^n \frac{\partial a_{jk}(y)}{\partial y_j} = \nabla_y \cdot (a(y)e_k) = \nabla_y \cdot (a(y)\nabla y_k). \quad (5.57)$$

This can be written as

$$-\nabla_y \cdot (a(y)\nabla_y(\chi_k + y_k)) = 0. \quad (5.58)$$

Multiplying by a test function  $\phi \in C^1(\mathbb{T}^n)$  and integrating by parts, using the fact that  $\phi$  and  $\nabla(\chi_k + y_k)$  are periodic (even though the function  $\chi_k + y_k$  is not periodic), gives

$$0 = \int_{\mathbb{T}^n} (a(y)\nabla_y(\chi_k + y_k) \cdot \nabla\phi)dy = 0,$$

as claimed.

Equation (5.58) gives a different representation for the effective diffusion matrix  $\bar{a}$ . Let us write

$$\begin{aligned}\bar{a}_{kj} &= \int_{\mathbb{T}^n} \left( a_{kj}(y) + \sum_{m=1}^n a_{km}(y) \frac{\partial \chi_j(y)}{\partial y_m} \right) dy = \int_{\mathbb{T}^n} (\nabla y_k \cdot a \nabla y_j + \nabla y_k \cdot a \nabla \chi_j) dy \\ &= \int_{\mathbb{T}^n} (\nabla y_k \cdot a \nabla (y_j + \chi_j)) dy = L(y_k, y_j + \chi_j).\end{aligned}\tag{5.59}$$

However, Proposition 5.2 implies, in particular, that

$$L(\chi_k, \chi_j + y_j) = 0.$$

It follows that

$$\bar{a}_{kj} = L(y_k + \chi_k, y_j + \chi_j).\tag{5.60}$$

This formula for  $\bar{a}$  shows that  $\bar{a}$  is a symmetric matrix since  $a(x)$  is symmetric.

Let us now show that the matrix  $\bar{a}_{ij}$  is positive-definite. Take an arbitrary vector  $\xi \in \mathbb{R}^n$ , then (5.60) implies that

$$(\bar{a}\xi \cdot \xi) = \int_{\mathbb{T}^n} (a \nabla (y_k + \chi_k) \cdot \nabla (y_j + \chi_j)) \xi_k \xi_j dy = \int_{\mathbb{T}^n} (a \nabla (\xi \cdot (y + \chi)) \cdot \nabla (\xi \cdot (y + \chi))) dy \geq 0.\tag{5.61}$$

In order to see that  $\bar{a}_{ij}$  is not only non-negative but is actually positive-definite, assume that  $(\bar{a}\xi \cdot \xi) = 0$ . Since the matrix  $a(y)$  is uniformly elliptic, it follows that

$$\nabla (\xi \cdot (y + \chi)) \equiv 0,$$

meaning that

$$\xi \cdot (y + \chi) = c\tag{5.62}$$

is a constant. As the function  $\chi(y)$  is periodic the only possibility for (5.62) to hold is that  $\xi \cdot y$  is uniformly bounded for all  $y \in \mathbb{T}^n$ , which means that  $\xi = 0$ .

As the matrix  $\bar{a}$  is positive definite, the homogenized problem (5.50) is well-posed.

### The two-scale convergence

An interesting way to formulate the homogenization theorem both for elliptic and other problems is via the two-scale convergence, a notion well suited for multiple scale problems.

**Definition 5.3** *Let  $u^\varepsilon$  be a family in  $L^2(\Omega)$ . We say that  $u^\varepsilon$  two-scale converges to  $u_0(x, y) \in L^2(\Omega \times \mathbb{T}^n)$  and write  $u^\varepsilon \xrightarrow{2} u$  if for any test function  $\phi \in L^2(\Omega; C(\mathbb{T}^n))$  we have*

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} u^\varepsilon(x) \phi(x, \frac{x}{\varepsilon}) dx = \int_{\Omega} \int_{\mathbb{T}^n} u(x, y) \phi(x, y) dy dx.\tag{5.63}$$

The reason we take the test functions in the space  $L^2(\Omega; C(\mathbb{T}^n))$  and not simply in  $L^2(\Omega \times \mathbb{T}^n)$  is that, in general, the restriction  $\phi(x, x/\varepsilon)$  is not necessarily a measurable function of  $x$  for a function  $\phi(x, y) \in L^2(\Omega \times \mathbb{T}^n)$ . These issues are discussed in detail in the fundamental paper by Gregoire Allaire “Homogenization and two-scale convergence”.

The idea behind the two-scale convergence can be seen on the following example. Consider the family

$$u^\varepsilon(x) = \psi(x) \sin\left(\frac{2\pi x}{\varepsilon}\right),$$

with  $\psi \in C_c^\infty$ . This family has no strong limit in  $L^2(\mathbb{R})$  and  $u^\varepsilon \rightarrow 0$  weakly in  $L^2(\mathbb{R})$ . The limit  $u(x) \equiv 0$  has no information whatsoever neither about the smooth component  $\psi(x)$  nor about the small scale oscillations  $\sin(2\pi x/\varepsilon)$ . As we will see, the two-scale convergence captures both.

On the other hand, it is important that the test functions in Definition 5.3 oscillate on exactly the same scale as the family  $u^\varepsilon(x)$ . Indeed, if we take a family  $u^\varepsilon(x) = u(x) \sin(\frac{2\pi x}{\varepsilon^2})$  then it is straightforward to check that for any test function  $\psi(x, y) \in L^2(\mathbb{R}; C(\mathbb{T}^n))$  we have

$$\int_{\mathbb{R}} u^\varepsilon(x) \psi\left(x, \frac{x}{\varepsilon}\right) dx \rightarrow 0 \text{ as } \varepsilon \rightarrow 0,$$

and similarly for the family  $v^\varepsilon(x) = v(x) \sin(\frac{2\pi x}{\sqrt{\varepsilon}})$  we have

$$\int_{\mathbb{R}} u^\varepsilon(x) \psi\left(x, \frac{x}{\varepsilon}\right) dx = \int_{\mathbb{R}} v(x) \psi\left(x, \frac{x}{\varepsilon}\right) \sin\left(\frac{2\pi x}{\sqrt{\varepsilon}}\right) dx \rightarrow 0.$$

Therefore, we need to take the test functions precisely on the scale of oscillations of  $u^\varepsilon(x)$  in order to see a non-trivial limit.

Two-scale convergence is stronger than the weak convergence in  $L^2(\Omega)$  as seen from the following.

**Proposition 5.4** *Let  $u^\varepsilon \in L^2(\Omega)$  be a two-scale convergent sequence:  $u^\varepsilon \xrightarrow{2} u \in L^2(\Omega \times \mathbb{T}^n)$ , then  $u^\varepsilon \rightarrow \bar{u}$  weakly in  $L^2(\Omega)$ , where*

$$\bar{u}(x) = \int_{\mathbb{T}^n} u(x, y) dy.$$

**Proof.** Take a test function  $\phi(x) \in L^2(\Omega)$  that is independent of  $y$ . Then we have

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} u^\varepsilon(x) \phi(x) dx = \int_{\Omega \times \mathbb{T}^n} u(x, y) \phi(x) dx dy = \int_{\Omega} \bar{u}(x) \phi(x) dx,$$

which implies that  $u^\varepsilon(x)$  converges weakly in  $L^2(\Omega)$  to  $\bar{u}(x)$ .

**Corollary 5.5** *Assume that a sequence  $u^\varepsilon(x) \in L^2(\Omega)$  two-scale converges to a limit  $u(x)$  that is independent of  $y$ . Then the weak  $L^2$  limit of  $u^\varepsilon(x)$  and the two-scale limit of  $u^\varepsilon(x)$  coincide.*

Here is how two-scale convergence is connected to multiple scale expansions, and to the discussion following Definition 5.3.

**Proposition 5.6** Consider a function  $u^\varepsilon \in C(\bar{\Omega})$  of the form

$$u^\varepsilon(x) = u_0(x, \frac{x}{\varepsilon}) + \varepsilon u_1(x, \frac{x}{\varepsilon}), \quad (5.64)$$

where  $u_{0,1} \in C(\bar{\Omega}; C(\mathbb{T}^n))$  and  $\Omega$  is a bounded open subset of  $\mathbb{R}^n$ . Then  $u_\varepsilon \xrightarrow{2} u_0$ .

**Proof.** Let  $\phi \in L^2(\Omega; C(\mathbb{T}^n))$  and define  $f_j(x, y) = u_j(x, y)\phi(x, y)$ ,  $j = 0, 1$ , as well as  $f_j^\varepsilon(x) = f_j(x, x/\varepsilon)$ . Then we have

$$\int_{\Omega} u^\varepsilon(x)\phi(x, \frac{x}{\varepsilon})dx = \int_{\Omega} f_0^\varepsilon(x) + \varepsilon \int_{\Omega} f_1^\varepsilon(x)dx. \quad (5.65)$$

The family  $f_0^\varepsilon$  converges weakly in  $L^2(\Omega)$  as  $\varepsilon \rightarrow 0$  to (this is a reasonably straightforward generalization of Theorem 5.1)

$$\bar{f}_0(x) = \int_{\mathbb{T}^n} f_0(x, y)dy,$$

hence for any function  $\psi \in L^2(\Omega)$  we have

$$\int_{\Omega} \psi(x)f_0^\varepsilon(x)dx \rightarrow \int_{\Omega \times \mathbb{T}^n} \psi(x)f_0(x, y)dy.$$

Taking  $\psi(x) \equiv 1$  gives

$$\int_{\Omega} f_0^\varepsilon(x)dx = \int_{\Omega} u_0(x, \frac{x}{\varepsilon})\phi(x, \frac{x}{\varepsilon})dx \rightarrow \int_{\Omega \times \mathbb{T}^n} f_0(x, y)dy = \int_{\Omega \times \mathbb{T}^n} u_0(x, y)\phi(x, y)dy.$$

On the other hand, for the second integral in the right side of (5.65) we note that, once again, the sequence  $f_1^\varepsilon$  is bounded in  $L^2(\Omega)$ , whence

$$\left| \int_{\Omega} f_1^\varepsilon(x)dx \right| \leq C\varepsilon \rightarrow 0, \text{ as } \varepsilon \rightarrow 0.$$

We conclude that

$$\int_{\Omega} u^\varepsilon(x)\phi(x, \frac{x}{\varepsilon})dx \rightarrow \int_{\Omega \times \mathbb{T}^n} u_0(x, y)\phi(x, y)dx dy, \quad (5.66)$$

finishing the proof.

Proposition 5.6 is very important – it says that if  $u^\varepsilon(x)$  indeed has only two scales of oscillations – the macroscopic scale of order  $O(1)$  and the microscopic scale  $O(\varepsilon)$  – then the two-scale limit is a much better suited notion than weak convergence in  $L^2(\Omega)$  since it retains the information about the small scale oscillations.

The next theorem provides a criterion for compactness.

**Theorem 5.7** Let  $u^\varepsilon$  be a bounded sequence in  $L^2(\Omega)$ . Then there exists a subsequence  $u^{\varepsilon_k}$  and a function  $u_0 \in L^2(\Omega \times \mathbb{T}^n)$  so that  $u^{\varepsilon_k} \xrightarrow{2} u_0$ .

**Proof.** We begin with the following lemma.

**Lemma 5.8** *Let  $\psi(x, y) \in L^2(\Omega; C(\mathbb{T}^n))$ , then  $\psi(x, x/\varepsilon)$  satisfies*

$$\|\psi(x, x/\varepsilon)\|_{L^2(\Omega)} \leq \|\psi(x, y)\|_{L^2(\Omega; C(\mathbb{T}^n))} = \left( \int_{\Omega} \sup_{y \in \mathbb{T}^n} |\psi(x, y)|^2 dx \right)^{1/2}, \quad (5.67)$$

and

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} |\psi(x, \frac{x}{\varepsilon})|^2 dx = \int_{\Omega \times \mathbb{T}^n} |\psi(x, y)|^2 dx dy. \quad (5.68)$$

Apart from the subtle measurability issues (5.67) is rather trivial, while (5.68) is a slightly stronger version of Theorem 5.1 – we refer to Gregoire Allaire’s paper “Homogenization and two-scale convergence” that we have already mentioned for the technicalities.

Let us now prove Theorem 5.7. Let  $u_\varepsilon$  be a bounded sequence in  $L^2(\Omega)$  with  $\|u^\varepsilon\|_{L^2(\Omega)} \leq C$ . According to Lemma 5.8, given any function  $\psi(x, y) \in L^2(\Omega; C(\mathbb{T}^n))$ , the functions  $\psi_\varepsilon(x) = \psi(x, x/\varepsilon)$  belong to  $L^2(\Omega)$ , and then

$$\left| \int_{\Omega} u^\varepsilon(x) \psi_\varepsilon(x) dx \right| \leq C \|\psi_\varepsilon\|_{L^2(\Omega)} \leq C \|\psi(x, y)\|_{L^2(\Omega; C(\mathbb{T}^n))}. \quad (5.69)$$

Therefore, for each  $\varepsilon > 0$  fixed,

$$A_\varepsilon(\psi) = \int_{\Omega} u^\varepsilon(x) \psi(x, \frac{x}{\varepsilon}) dx$$

is a bounded linear functional on  $L^2(\Omega; C(\mathbb{T}^n))$ . The dual space of  $L^2(\Omega; C(\mathbb{T}^n))$  is  $L^2(\Omega; M(\mathbb{T}^n))$  where  $M(\mathbb{T}^n)$  is the space of bounded signed Radon measures on  $\mathbb{T}^n$ . Therefore, there exists a bounded function  $\mu_\varepsilon \in L^2(\Omega; M(\mathbb{T}^n))$  so that

$$\int_{\Omega} \psi(x, y) \mu_\varepsilon(x, dy) dx = \int_{\Omega} u_\varepsilon(x) \psi_\varepsilon(x) dx.$$

Hence, there exists a sub-sequence  $\varepsilon_k \rightarrow 0$  so that  $\mu_{\varepsilon_k} \rightarrow \mu_0$  in  $L^2(\Omega; M(\mathbb{T}^n))$  in the weak\* topology, that is, for any  $\psi \in L^2(\Omega; C(\mathbb{T}^n))$  we have

$$\int_{\Omega} \psi(x, y) \mu_{\varepsilon_k}(x, dy) dx \rightarrow \int_{\Omega} \psi(x, y) \mu_0(x, dy) dx,$$

as  $\varepsilon_k \rightarrow 0$ . In other words, we have

$$\int_{\Omega} u_\varepsilon(x) \psi_\varepsilon(x) dx \rightarrow \int_{\Omega} \psi(x, y) \mu_0(x, dy) dx. \quad (5.70)$$

We would be done if we knew that the measure  $\mu(x, dy)$  has the form

$$\mu(x, dy) = u_0(x, y) dx dy,$$

as then (5.70) would simply say that  $u_\varepsilon(x)$  two-scale converge to  $u_0(x)$ . We now show that, indeed,  $\mu_0(x, dy)$  has a density. On the other hand, we also know that

$$\|\psi_\varepsilon(x)\|_{L^2(\Omega)}^2 = \int_{\Omega} |\psi(x, \frac{x}{\varepsilon})|^2 dx \rightarrow \int_{\Omega \times \mathbb{T}^n} |\psi(x, y)|^2 dx dy.$$

Therefore, we may pass to the limit in the first inequality in (5.69) to get

$$\left| \int_{\Omega \times \mathbb{T}^n} \psi(x, y) \mu_0(x, dy) dx \right| \leq C \int_{\Omega \times \mathbb{T}^n} |\psi(x, y)|^2 dx dy,$$

for any  $\psi \in L^2(\Omega; C(\mathbb{T}^n))$ . Therefore, the Riesz Representation theorem implies that measure  $\mu_0(dx, y)$  does have density

$$u_0(x, y) \in L^2(\Omega \times \mathbb{T}^n),$$

that is,

$$\int_{\Omega} \psi(x, y) \mu_0(x, dy) dx = \int_{\Omega \times \mathbb{T}^n} u_0(x, y) \psi(x, y) dx dy.$$

Summarizing, we have shown that for any function  $\psi(x, y) \in L^2(\Omega; C(\mathbb{T}^n))$  we have

$$\int_{\Omega} u^\varepsilon(x) \psi(x, \frac{x}{\varepsilon}) dx \rightarrow \int_{\Omega \times \mathbb{T}^n} u_0(x, y) \psi(x, y) dx dy,$$

whence  $u^\varepsilon \xrightarrow{2} u_0$ , and the proof is complete.

The next theorem shows that two-scale convergence is “strong” under an additional assumption.

**Theorem 5.9** *Let  $u_\varepsilon \in L^2(\Omega)$  two-scale converge to  $u_0(x, y) \in L^2(\Omega \times \mathbb{T}^n)$ . Assume, in addition, that*

$$\lim_{\varepsilon \rightarrow 0} \|u_\varepsilon\|_{L^2(\Omega)} = \|u_0\|_{L^2(\Omega \times \mathbb{T}^n)}. \quad (5.71)$$

*Then for any sequence  $v_\varepsilon(x)$  that two-scale converges to a limit  $v_0(x, y) \in L^2(\Omega \times \mathbb{T}^n)$  we have*

$$u_\varepsilon(x) v_\varepsilon(x) \rightarrow \int_{\mathbb{T}^n} u_0(x, y) v_0(x, y) dy, \quad (5.72)$$

*in the sense of distributions on  $\Omega$ . Moreover, if  $u_0(x, y) \in L^2(\Omega; C(\mathbb{T}^n))$  then*

$$\lim_{\varepsilon \rightarrow 0} \left\| u_\varepsilon(x) - u_0(x, \frac{x}{\varepsilon}) \right\|_{L^2(\Omega)} = 0. \quad (5.73)$$

Note that condition (5.71) is not very restrictive: for instance, it holds for functions of the form  $u_\varepsilon(x) = \psi(x) \eta(x/\varepsilon)$ , and more generally  $u_\varepsilon(x) = \psi(x, x/\varepsilon)$  with  $\psi(x, y) \in L^2(\Omega; C(\mathbb{T}^n))$ . Physically, this condition means that  $u_\varepsilon(x)$  oscillate exactly on the scale  $\varepsilon$ , and not on scales much larger or much smaller than  $\varepsilon$  – if that were the case we would have lost mass in the limit, as we have seen before.

Another crucial aspect of this theorem is it allows to pass to the limit in the product of two-scale convergent sequences (if condition (5.71) holds for one of them) – this is drastically different from the weak convergence in  $L^2(\Omega)$  – the product of two weakly convergent sequences need not converge to the product of the limits.

**Proof.** Let  $\psi_n(x, y)$  be a sequence of smooth functions in  $L^2(\Omega; C(\mathbb{T}^n))$  that converge strongly in  $L^2(\Omega \times \mathbb{T}^n)$  to  $u_0(x, y)$ . By definition of two-scale convergence and using (5.71)

we have:

$$\begin{aligned}
\lim_{\varepsilon \rightarrow 0} \int_{\Omega} |u^{\varepsilon}(x) - \psi_n(x, \frac{x}{\varepsilon})|^2 dx &= \lim_{\varepsilon \rightarrow 0} \int_{\Omega} (|u^{\varepsilon}(x)|^2 - 2u^{\varepsilon}(x)\psi_n(x, \frac{x}{\varepsilon}) + |\psi_n(x, \frac{x}{\varepsilon})|^2) dx \\
&= \int_{\Omega \times \mathbb{T}^n} (|u_0(x, y)|^2 - 2u_0(x, y)\psi_n(x, y) + |\psi_n(x, y)|^2) dx dy \\
&= \int_{\Omega \times \mathbb{T}^n} |u_0(x, y) - \psi_n(x, y)|^2 dx dy.
\end{aligned} \tag{5.74}$$

In particular, if  $u_0(x, y)$  itself is in  $L^2(\Omega; C(\mathbb{T}^n))$  then this is nothing but (5.73).

Next, if  $v_{\varepsilon}$  is another two-scale convergent sequence so that  $v_{\varepsilon} \xrightarrow{2} v_0$  then for any smooth test function  $\phi(x, y)$  we have

$$\int_{\Omega} \phi(x) u_{\varepsilon}(x) v_{\varepsilon}(x) dx = \int_{\Omega} \phi(x) \psi_n(x, \frac{x}{\varepsilon}) v_{\varepsilon}(x) dx + \int_{\Omega} \phi(x) \left[ u_{\varepsilon}(x) - \psi_n(x, \frac{x}{\varepsilon}) \right] v_{\varepsilon}(x) dx.$$

We may now pass to the limit  $\varepsilon \rightarrow 0$  taking into account (5.74) and recalling that the sequence  $v_{\varepsilon}(x)$  is uniformly bounded in  $L^2(\Omega)$  as implied by the two-scale convergence:

$$\begin{aligned}
&\left| \lim_{\varepsilon \rightarrow 0} \left[ \int_{\Omega} \phi(x) u_{\varepsilon}(x) v_{\varepsilon}(x) dx - \int_{\Omega} \phi(x) \psi_n(x, \frac{x}{\varepsilon}) v_{\varepsilon}(x) dx \right] \right| \\
&\leq \lim_{\varepsilon \rightarrow 0} \int_{\Omega} |\phi(x)| \left| u_{\varepsilon}(x) - \psi_n(x, \frac{x}{\varepsilon}) \right| |v_{\varepsilon}(x)| dx \leq C \lim_{\varepsilon \rightarrow 0} \|u_{\varepsilon}(x) - \psi_n(x, \frac{x}{\varepsilon})\|_{L^2(\Omega)}.
\end{aligned}$$

We now pass to the limit  $n \rightarrow +\infty$  with the help of (5.74) and conclude that

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} \phi(x) u_{\varepsilon}(x) v_{\varepsilon}(x) dx = \int_{\Omega \times \mathbb{T}^n} \phi(x) u_0(x, y) v_0(x, y) dy,$$

which is exactly (5.72).

The next theorem accounts for what happens not just with the functions but with their derivatives under the two-scale convergence.

**Theorem 5.10** (i) *Let  $u_{\varepsilon}$  be a bounded sequence in  $H^1(\Omega)$  that converges weakly to a limit  $u$  in  $H^1(\Omega)$ . Then  $u_{\varepsilon}$  two-scale converges to  $u(x)$  and there exists a function  $u_1(x, y)$  in  $L^2(\Omega; H^1(\mathbb{T}^n))$  such that, up to extraction of a subsequence,  $\nabla u_{\varepsilon}$  two-scale converges to  $\nabla_x u(x) + \nabla_y u_1(x, y)$ .*

(ii) *Let  $u_{\varepsilon}$  and  $\varepsilon \nabla u_{\varepsilon}$  be both bounded in  $L^2(\Omega)$ . Then there exists a function  $u_0(x, y) \in L^2(\Omega; H^1(\mathbb{T}^n))$  so that, up to extraction of a subsequence,  $u_{\varepsilon}$  and  $\varepsilon \nabla u_{\varepsilon}$  two-scale converge to  $u_0(x, y)$  and  $\nabla_y u_0(x, y)$ , respectively.*

The first part of the theorem is perfectly suited for functions that can be represented as

$$u_{\varepsilon}(x) = u_0(x) + \varepsilon u_1(x, \frac{x}{\varepsilon}) + O(\varepsilon^2),$$

then, provided we have bounds on the remainder,  $u_{\varepsilon}$  two-scale converges to  $u_0(x)$  and  $\nabla u_{\varepsilon}$  two-scale converges to  $\nabla_x u_0(x) + \nabla_y u_1(x, x/\varepsilon)$ . The second part is suited for functions of the form

$$u_{\varepsilon}(x) = u_0(x, \frac{x}{\varepsilon}) + \varepsilon u_1(x, \frac{x}{\varepsilon}) + O(\varepsilon^2),$$

where the leading term is also oscillating on the scale  $\varepsilon$ .

**Proof.** (i) Since  $u_\varepsilon$  and  $\nabla u_\varepsilon$  are bounded in  $L^2(\Omega)$ , up to extraction of a subsequence, they converge to the limits  $u_0 \in L^2(\Omega \times \mathbb{T}^n)$  and  $\chi_0 \in (L^2(\Omega \times \mathbb{T}^n))^n$ . Thus, for any scalar smooth test function  $\phi(x, y)$  and vector-valued smooth test function  $\eta(x, y)$  we have

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} u_\varepsilon(x) \phi(x, \frac{x}{\varepsilon}) dx = \int_{\Omega \times \mathbb{T}^n} u_0(x, y) \phi(x, y) dx dy, \quad (5.75)$$

and

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} \nabla u_\varepsilon(x) \cdot \eta(x, \frac{x}{\varepsilon}) dx = \int_{\Omega \times \mathbb{T}^n} \chi_0(x, y) \cdot \eta(x, y) dx dy. \quad (5.76)$$

Integrating by parts, on the other hand, gives

$$\varepsilon \int_{\Omega} \nabla u_\varepsilon(x) \cdot \eta(x, \frac{x}{\varepsilon}) dx = - \int_{\Omega} u_\varepsilon(x) \left[ \nabla_y \cdot \eta(x, \frac{x}{\varepsilon}) + \varepsilon \nabla_x \cdot \eta(x, \frac{x}{\varepsilon}) \right] dx.$$

Letting  $\varepsilon \rightarrow 0$  and using (5.75) leads to

$$0 = - \int_{\Omega \times \mathbb{T}^n} u_0(x, y) \nabla_y \cdot \eta(x, y) dx dy,$$

for all vector-valued test functions  $\eta(x, y)$ . It follows that  $u_0(x, y)$  does not depend on  $x$ . Moreover, the weak convergence of  $u_\varepsilon$  to  $u$  implies that

$$u(x) = \int_{\mathbb{T}^n} u_0(x, y) dy,$$

and since  $u_0$  does not depend on  $y$  we conclude that  $u(x) = u_0(x)$ , proving the first claim in (i). In order to show that there exists a function  $u_1(x, y) \in H^1(\Omega; H^1(\mathbb{T}^n))$  such that  $\nabla u_\varepsilon$  two-scale converges to  $\nabla_x u(x) + \nabla_y u_1(x, y)$  we choose  $\eta(x, y)$  in (5.76) such that  $\nabla_y \cdot \eta(x, y) = 0$ :

$$\begin{aligned} \int_{\Omega \times \mathbb{T}^n} \chi_0(x, y) \cdot \eta(x, y) dx dy &= \lim_{\varepsilon \rightarrow 0} \int_{\Omega} \nabla u_\varepsilon(x) \cdot \eta(x, \frac{x}{\varepsilon}) dx = - \lim_{\varepsilon \rightarrow 0} \int_{\Omega} u_\varepsilon(x) \nabla_x \cdot \eta(x, \frac{x}{\varepsilon}) dx \\ &= - \int_{\Omega \times \mathbb{T}^n} u(x) \nabla_x \cdot \eta(x, y) dx dy = \int_{\Omega \times \mathbb{T}^n} \nabla u(x) \cdot \eta(x, y) dx dy. \end{aligned} \quad (5.77)$$

Therefore, for any vector-valued smooth test function  $\eta(x, y)$  such that  $\nabla_y \cdot \eta(x, y) = 0$  we have

$$\int_{\Omega \times \mathbb{T}^n} [\chi_0(x, y) - \nabla u(x)] \cdot \eta(x, y) dx dy = 0. \quad (5.78)$$

This means that there exists some function  $u_1(x, y) \in L^2(\Omega; H^1(\mathbb{T}^n))$  such that

$$\chi_0(x, y) - \nabla u(x) = \nabla_y u_1(x, y),$$

which is exactly the second claim in part (i).

We now prove part (ii) of the theorem. Since  $u_\varepsilon$  and  $\varepsilon\nabla u_\varepsilon$  are bounded in  $L^2(\Omega)$ , up to extraction of a subsequence, they converge to the limits  $u_0 \in L^2(\Omega \times \mathbb{T}^n)$  and  $\chi_0 \in (L^2(\Omega \times \mathbb{T}^n))^n$ . Thus, for any scalar smooth test function  $\phi(x, y)$  and vector-valued smooth test function  $\eta(x, y)$  we have

$$\lim_{\varepsilon \rightarrow 0} \int_{\Omega} u_\varepsilon(x) \phi(x, \frac{x}{\varepsilon}) dx = \int_{\Omega \times \mathbb{T}^n} u_0(x, y) \phi(x, y) dx dy, \quad (5.79)$$

and

$$\lim_{\varepsilon \rightarrow 0} \varepsilon \int_{\Omega} \nabla u_\varepsilon(x) \cdot \eta(x, \frac{x}{\varepsilon}) dx = \int_{\Omega \times \mathbb{T}^n} \chi_0(x, y) \cdot \eta(x, y) dx dy. \quad (5.80)$$

Integrating by parts, on the other hand, gives

$$\varepsilon \int_{\Omega} \nabla u_\varepsilon(x) \cdot \eta(x, \frac{x}{\varepsilon}) dx = - \int_{\Omega} u_\varepsilon(x) \left[ \nabla_y \cdot \eta(x, \frac{x}{\varepsilon}) + \varepsilon \nabla_x \cdot \eta(x, \frac{x}{\varepsilon}) \right] dx,$$

which after passing to the limit  $\varepsilon \rightarrow 0$  implies that

$$\int_{\Omega} \chi_0(x, y) \cdot \eta(x, y) dx dy = - \int_{\Omega} u_0(x, y) \nabla_y \cdot \eta(x, y) dx dy,$$

for all smooth vector-valued test functions  $\eta(x, y)$ . It follows that  $\chi_0(x, y) = \nabla_y u_0(x, y)$  and the proof of part (ii) of the theorem is complete.

## Two-scale convergence in elliptic homogenization

As we have seen in the preceding analysis, we expect, at least on the level of a formal asymptotic analysis, that the solution of the boundary value problem

$$\begin{aligned} -\nabla \cdot \left( a \left( \frac{x}{\varepsilon} \right) \nabla \phi^\varepsilon \right) &= f \text{ in } \Omega, \\ \phi^\varepsilon &= 0 \text{ on } \partial\Omega, \end{aligned} \quad (5.81)$$

is well approximated by the solution of the homogenized problem

$$\begin{aligned} -\nabla \cdot (\bar{a} \nabla \bar{\phi}) &= f \text{ in } \Omega, \\ \bar{\phi} &= 0 \text{ on } \partial\Omega, \end{aligned} \quad (5.82)$$

with the effective diffusion matrix  $\bar{a}$  given by (5.59). Moreover, we expect that  $\phi^\varepsilon(x)$  can be approximated by a multiple scales expansion

$$\phi^\varepsilon(x) = \bar{\phi}(x) + \varepsilon \phi_1^\varepsilon \left( x, \frac{x}{\varepsilon} \right) + \varepsilon^2 \phi_2 \left( x, \frac{x}{\varepsilon} \right) + \dots \quad (5.83)$$

with the function

$$\phi_1(x, y) = \sum_{j=1}^n \chi_j(y) \frac{\partial \bar{\phi}(x)}{\partial x_j} \quad (5.84)$$

that is periodic in the fast variable  $y \in \mathbb{T}^n$  and is reasonably nice in the slow variable  $x \in \Omega$ . The notion of the two-scale convergence allows us to formalize this statement.

Let us define the space

$$H = \left\{ u \in H^1(\mathbb{T}^n) : \int_{\mathbb{T}^n} u(y) dy = 0 \right\},$$

as well as

$$X = H_0^1(\Omega) \times L^2(\Omega; H),$$

so that an element of  $X$  is  $(u(x), v(x, y))$  with the function  $u \in H_0^1(\Omega)$  and the function  $v(x, y)$  that is periodic in  $y$  and satisfies

$$\int_{\Omega \times \mathbb{T}^n} |\nabla_y v(x, y)|^2 dx dy < +\infty.$$

This is a Hilbert space with the inner product

$$\langle U, V \rangle_X = \int_{\Omega} (\nabla_x u(x) \cdot \nabla_x v(x)) dx + \int_{\Omega \times \mathbb{T}^n} (\nabla_y u_1(x, y) \cdot \nabla_y v_1(x, y)) dx dy,$$

for all  $U = (u, u_1)$ ,  $V = (v, v_1)$ . The corresponding norm is

$$\|U\|_X^2 = \|\nabla u\|_{L^2(\Omega)}^2 + \|\nabla_y u_1\|_{L^2(\Omega \times \mathbb{T}^n)}^2.$$

Our first result is that the solution of (5.81) has a two-scale limit.

**Proposition 5.11** *Let  $\phi^\varepsilon(x)$  be the solution of (5.81), then there exist a pair of functions  $(\phi(x), \phi_1(x, y)) \in X$  and a subsequence  $\varepsilon_k \rightarrow 0$  so that  $\phi^\varepsilon$  and  $\nabla \phi^\varepsilon(x)$  two-scale converge, along a subsequence to  $u(x)$  and  $\nabla_x u + \nabla_y u_1$ ,*

**Proof.** Multiplying

$$-\nabla \cdot \left( a\left(\frac{x}{\varepsilon}\right) \phi^\varepsilon \right) = f,$$

by  $\phi^\varepsilon$  and integrating by parts using the boundary condition  $\phi^\varepsilon = 0$  on  $\partial\Omega$  we get

$$\int_{\Omega} a\left(\frac{x}{\varepsilon}\right) \nabla \phi^\varepsilon \cdot \nabla \phi^\varepsilon dx = \int_{\Omega} f \phi^\varepsilon dx \leq \|f\|_{L^2} \|\phi^\varepsilon\|_{L^2}.$$

As the matrix  $a(y)$  is uniformly positive-definite we deduce that

$$\int_{\Omega} |\nabla \phi^\varepsilon|^2 dx \leq c_0 \|f\|_{L^2} \|\phi^\varepsilon\|_{L^2}. \quad (5.85)$$

As  $\phi^\varepsilon(x) = 0$  on the boundary, it satisfies the Poincaré inequality

$$\int_{\Omega} |\phi^\varepsilon(x)|^2 dx \leq C_p \int_{\Omega} |\nabla \phi^\varepsilon(x)|^2 dx, \quad (5.86)$$

with the constant  $C_p$  that depends only on the domain  $\Omega$ . Using this in (5.85) gives

$$\int_{\Omega} |\nabla \phi^\varepsilon|^2 dx \leq c_0 C_p \|f\|_{L^2} \|\nabla \phi^\varepsilon\|_{L^2}, \quad (5.87)$$

whence

$$\|\nabla\phi_\varepsilon\|_{L^2} \leq c_0 C_p \|f\|_{L^2}. \quad (5.88)$$

Hence, the sequence  $\phi^\varepsilon(x)$  is uniformly bounded in  $H_0^1(\Omega)$ , so that we may apply part (i) of Theorem 5.10 to finish the proof of this lemma.

We now show that the limits  $u(x)$  and  $u_1(x, y)$  satisfy the two-scale system

$$\begin{aligned} -\nabla_y \cdot (a(y)(\nabla_x u + \nabla_y u_1)) &= 0 \text{ in } \Omega \times \mathbb{T}^n, \\ -\nabla_x \cdot \left( \int_{\mathbb{T}^n} a(y)(\nabla_x u + \nabla_y u_1) dy \right) &= f \text{ in } \Omega, \end{aligned} \quad (5.89)$$

with the boundary condition

$$u(x) = 0 \text{ for } x \in \partial\Omega, \quad u_1(x, y) \text{ is periodic in } y. \quad (5.90)$$

This is done nicely via the weak formulation. Consider the bilinear form, defined for  $U = (u, u_1)$  and  $\Phi = (\phi_0, \phi_1) \in X$ :

$$\mathcal{L}[U, \Phi] = \int_{\Omega \times \mathbb{T}^n} (a(y)(\nabla_x u + \nabla_y u_1) \cdot (\nabla_x \phi_0 + \nabla_y \phi_1)) dy dx.$$

The weak form of (5.89) is to find  $U = (u, u_1) \in X$  such that

$$\mathcal{L}[U, \Phi] = \langle f, \phi_0 \rangle, \quad \text{all } \Phi \in X. \quad (5.91)$$

To see that, first take  $\Phi$  with  $\phi_0 = 0$ , then (5.91) says that

$$\int_{\Omega \times \mathbb{T}^n} (a(y)(\nabla_x u + \nabla_y u_1) \cdot \nabla_y \phi_1) dy dx = 0, \quad (5.92)$$

which is the weak form of the first equation in (5.89). Next, taking  $\Phi$  with  $\phi_1 = 0$  gives

$$\int_{\Omega \times \mathbb{T}^n} (a(y)(\nabla_x u + \nabla_y u_1) \cdot \nabla_x \phi_0) dy dx = \int_{\Omega} f(x) \phi_0(x) dx, \quad (5.93)$$

which is the weak form of the second equation in (5.89). The boundary condition for  $u(x)$  is incorporated in the definition of the space  $X$  – it requires that  $u \in H_0^1(\Omega)$ .

**Lemma 5.12** *Let  $u^\varepsilon(x)$  be a weak solution of (5.63). Then any limit point  $(u, u_1)$  from Proposition 5.11 is a weak solution of the two-scale system (5.89)-(5.90).*

**Proof.** The weak form of (5.63) is to find  $\phi^\varepsilon \in H_0^1(\Omega)$  such that

$$\int_{\Omega} (a(\frac{x}{\varepsilon}) \nabla \phi^\varepsilon \cdot \nabla \psi) dx = \langle f, \psi \rangle, \quad \text{for any } \psi \in H_0^1(\Omega). \quad (5.94)$$

Let us choose the test function  $\psi$  to be of the form

$$\psi(x) = \psi_0(x) + \varepsilon \psi_1(x, \frac{x}{\varepsilon}). \quad (5.95)$$

This is a very important idea and is close to what is known as the perturbed test function method – in that method you modify the test function so as to cancel potentially large terms. Here, we modify the test function to study the two-scale convergence. Inserting this test function into (5.94) gives

$$I_1 + \varepsilon I_2 = \langle f, \psi_0 + \varepsilon \psi_1^\varepsilon \rangle,$$

where  $\psi_1^\varepsilon(x) = \psi_1(x, x/\varepsilon)$ ,

$$I_1 = \int_{\Omega} (\nabla \phi^\varepsilon \cdot a(\frac{x}{\varepsilon})(\nabla_x \psi_0(x) + \nabla_y \psi_1(x, \frac{x}{\varepsilon}))) dx,$$

and

$$I_2 = \int_{\Omega} (\nabla \phi^\varepsilon \cdot a(\frac{x}{\varepsilon}) \nabla_x \psi_1(x, \frac{x}{\varepsilon})) dx.$$

Now, we recall that Proposition 5.11 ensures that  $\phi^\varepsilon(x)$  two-scale converges to  $u(x)$  and  $\nabla \phi^\varepsilon(x)$  converges to  $\nabla_x u(x) + \nabla_y u_1(x, y)$ . We can pass to the two-scale limit in the expressions for  $I_1$  and  $I_2$ , to obtain

$$I_1 \rightarrow \int_{\Omega} \int_{\mathbb{T}^n} (a(y)(\nabla_x u + \nabla_y u_1) \cdot (\nabla_x \psi_0 + \nabla_y \psi_1)) dy dx, \text{ as } \varepsilon \rightarrow 0,$$

and  $\varepsilon I_2 \rightarrow 0$ . Moreover, we have

$$\langle f, \psi_0 + \varepsilon \psi_1^\varepsilon \rangle \rightarrow \langle f, \psi_0 \rangle.$$

Putting these together gives

$$\int_{\Omega} \int_{\mathbb{T}^n} (a(y)(\nabla_x u + \nabla_y u_1) \cdot (\nabla_x \psi_0 + \nabla_y \psi_1)) dy dx = \langle f, \psi_0 \rangle,$$

which is, indeed, the weak form (5.91) of the two-scale system (5.89)-(5.90).

The next step is to show that the two-scale system has a unique solution.

**Lemma 5.13** *The two-scale system (5.89)-(5.90) has a unique solution  $(u, u_1) \in X$ .*

**Proof.** The proof is based on the Lax-Milgram theorem. In order to apply it to the weak formulation (5.91)

$$\mathcal{L}[U, \Phi] = \langle f, \phi_0 \rangle, \tag{5.96}$$

of the two-scale system we need to verify that the bi-linear form  $\mathcal{L}$  is coercive and continuous, that is, there exists a constant  $c_0 > 0$  so that

$$c_0 \|U\|_X \leq \mathcal{L}[U, U] \leq c_0^{-1} \|U\|_X. \tag{5.97}$$

Note that if  $U = (u(x), v(x, y))$  then

$$\|U\|_X^2 = \int_{\Omega \times \mathbb{T}^n} |\nabla_x u(x) + \nabla_y v(x, y)|^2 dx dy.$$

This is because

$$\begin{aligned}\|U\|_X^2 &= \int_{\Omega} |\nabla_x u(x)|^2 dx + \int_{\Omega \times \mathbb{T}^n} |\nabla_y v(x, y)|^2 dx dy = \int_{\Omega \times \mathbb{T}^n} (|\nabla_x u(x)|^2 + |\nabla_y v(x, y)|^2) dx dy \\ &= \int_{\Omega \times \mathbb{T}^n} |\nabla_x u(x) + \nabla_y v(x, y)|^2 dx dy,\end{aligned}$$

as

$$\int_{\Omega \times \mathbb{T}^n} \nabla_x u(x) \cdot \nabla_y v(x, y) dx dy = 0,$$

since the function  $v(x, y)$  is periodic in  $y$ . Now, in order to prove the upper bound in (5.97) we notice that

$$\begin{aligned}\langle \mathcal{L}U, U \rangle &= \int_{\Omega \times \mathbb{T}^n} (a(y)(\nabla_x u + \nabla_y v) \cdot (\nabla_x u + \nabla_y v)) dy dx \\ &\leq \|a\|_{L^\infty} \int_{\Omega \times \mathbb{T}^n} |\nabla_x u + \nabla_y v|^2 dx dy = C_0 \|U\|_X^2,\end{aligned}$$

because  $a(y)$  is uniformly bounded from above while uniform ellipticity of  $a(y)$  means that

$$\begin{aligned}\langle \mathcal{L}U, U \rangle &= \int_{\Omega \times \mathbb{T}^n} (a(y)(\nabla_x u + \nabla_y v) \cdot (\nabla_x u + \nabla_y v)) dy dx \\ &\geq c_0 \int_{\Omega \times \mathbb{T}^n} |\nabla_x u + \nabla_y v|^2 dx dy = c_0 \|U\|_X^2.\end{aligned}$$

Therefore, the bilinear form  $\mathcal{L}(U, \Phi)$  is, indeed, continuous and coercive, hence the Lax-Milgram lemma implies that the weak equation

$$\mathcal{L}[U, \Phi] = \langle f, \phi_0 \rangle, \text{ for all } \Phi = (\phi_0, \phi_1) \in X, \quad (5.98)$$

has a unique solution  $U \in X$ .

Next, we relate the two-scale system to the homogenized equation (5.82).

**Lemma 5.14** *Let  $(u, u_1) \in X$  be the unique solution of the two-scale system (5.89)-(5.90). Then  $u$  is the unique solution of the homogenized problem*

$$\begin{aligned}-\nabla \cdot (\bar{a} \nabla u) &= f \text{ in } \Omega, \\ u &= 0 \text{ on } \partial\Omega,\end{aligned} \quad (5.99)$$

and

$$u_1(x, y) = \chi(y) \cdot \nabla_x u(x), \quad (5.100)$$

where  $\chi(y)$  is the mean-zero solution of the cell problem.

**Proof.** Inserting expression (5.100) for  $u_1$  into (5.89) gives

$$-\frac{\partial}{\partial y_i} \left( a_{ij}(y) \left( \frac{\partial u}{\partial x_j} + \frac{\partial}{\partial y_j} (\chi_k(y) \frac{\partial u}{\partial x_k}) \right) \right) = 0 \text{ in } \Omega \times \mathbb{T}^n, \quad (5.101)$$

or

$$-\frac{\partial}{\partial y_i} \left( a_{ij}(y) \left( \frac{\partial \chi_k(y)}{\partial y_j} \right) \right) \frac{\partial u}{\partial x_k} = \frac{\partial a_{ik}(y)}{\partial y_i} \frac{\partial u}{\partial x_k} \text{ in } \Omega \times \mathbb{T}^n, \quad (5.102)$$

which holds because

$$-\frac{\partial}{\partial y_i} \left( a_{ij}(y) \frac{\partial \chi_k}{\partial y_j} \right) = \frac{\partial a_{ik}}{\partial y_i}.$$

On the other hand, (5.90) is

$$-\frac{\partial}{\partial x_j} \left( \int_{\mathbb{T}^n} a_{ij}(y) \left( \frac{\partial u}{\partial x_j} + \frac{\partial}{\partial y_j} (\chi_k(y) \frac{\partial u}{\partial x_k}) \right) dy \right) = f \text{ in } \Omega, \quad (5.103)$$

which is

$$-\frac{\partial}{\partial x_i} \left( \bar{a}_{ij} \frac{\partial u}{\partial x_j} \right) = f,$$

with

$$\bar{a}_{ij} = \int_{\mathbb{T}^n} (a_{ij}(y) + a_{im}(y) \frac{\partial \chi_j(y)}{\partial y_m}) dy. \quad (5.104)$$

Therefore, we have proved the following theorem.

**Theorem 5.15** *Let  $\phi^\varepsilon(x)$  be the solution of*

$$\begin{aligned} -\nabla \cdot \left( a \left( \frac{x}{\varepsilon} \right) \nabla \phi^\varepsilon \right) &= f, \quad x \in \Omega, \\ \phi^\varepsilon &= 0 \text{ on } \partial\Omega, \end{aligned} \quad (5.105)$$

*in a smooth domain  $\Omega$  and with  $f \in L^2(\Omega)$ . Let the matrix  $\bar{a}_{ij}$  be given by (5.104) with  $\chi_k(x)$ ,  $k = 1, \dots, n$  being the mean-zero periodic solution of*

$$-\nabla \cdot (a(y) \nabla \chi_k) = \sum_{j=1}^n \frac{\partial a_{jk}(y)}{\partial y_j}. \quad (5.106)$$

*Then  $u_\varepsilon(x)$  converges strongly in  $L^2(\Omega)$  and weakly in  $H_0^1(\Omega)$  to  $u(x)$ , solution of*

$$\begin{aligned} -\nabla \cdot (\bar{a} \nabla \phi^\varepsilon) &= f, \quad x \in \Omega, \\ \phi^\varepsilon &= 0 \text{ on } \partial\Omega. \end{aligned} \quad (5.107)$$

### Strong convergence in $H^1(\Omega)$

We now show how the  $L^2$ -convergence can be improved to  $H^1$ -convergence.

**Theorem 5.16** *Let  $\phi^\varepsilon(x)$  and  $u(x)$  be as in Theorem 5.15, and  $\Omega$  be a smooth domain, then*

$$\left\| \phi_\varepsilon(x) - \left( u(x) + \varepsilon \chi \left( \frac{x}{\varepsilon} \right) \cdot \nabla_x u(x) \right) \right\|_{H^1(\Omega)} \rightarrow 0, \quad (5.108)$$

*as  $\varepsilon \rightarrow 0$ .*

**Proof.** We need to prove that

$$\left\| \nabla \phi_\varepsilon(x) - \left( \nabla u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}) + \varepsilon \nabla_x u_1(x, \frac{x}{\varepsilon}) \right) \right\|_{L^2(\Omega)} \rightarrow 0, \quad (5.109)$$

where

$$u_1(x, y) = \chi(y) \cdot \nabla_x u(x).$$

As  $\|\varepsilon \nabla_x u_1(x, x/\varepsilon)\|_{L^2} \rightarrow 0$  as  $\varepsilon \rightarrow 0$ , it suffices to prove that

$$\left\| \nabla \phi_\varepsilon(x) - \left( \nabla u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}) \right) \right\|_{L^2(\Omega)} \rightarrow 0, \quad (5.110)$$

To this end, we proceed as follows:

$$\begin{aligned} & \int_{\Omega} \left| \nabla \phi_\varepsilon(x) - \left( \nabla u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}) \right) \right|^2 dx \\ & \leq C \int_{\Omega} a(\frac{x}{\varepsilon}) \left( \nabla \phi_\varepsilon(x) - \left( \nabla u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}) \right) \right) \cdot \left( \nabla \phi_\varepsilon(x) - \left( \nabla u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}) \right) \right) dx \\ & = C \int_{\Omega} (a(\frac{x}{\varepsilon}) \nabla \phi^\varepsilon \cdot \nabla \phi^\varepsilon) dx + C \int_{\Omega} (a(\frac{x}{\varepsilon}) (\nabla_x u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon})) \cdot (\nabla_x u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}))) dx \\ & \quad - 2C \int_{\Omega} (a(\frac{x}{\varepsilon}) \nabla \phi^\varepsilon \cdot (\nabla_x u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}))) dx = C \langle f, \phi^\varepsilon \rangle + I_1^\varepsilon + I_2^\varepsilon. \end{aligned}$$

We already know that  $\phi^\varepsilon$  converges strongly in  $L^2(\Omega)$  to  $u(x)$ , hence

$$\langle f, \phi^\varepsilon \rangle \rightarrow \langle f, u \rangle = \mathcal{L}[U, U],$$

where  $U = (u, u_1)$ . The last equality above follows from the two-scale system for  $U$ . Moreover, we have

$$\begin{aligned} & \int_{\Omega} (a(\frac{x}{\varepsilon}) (\nabla_x u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon})) \cdot (\nabla_x u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}))) dx \\ & \rightarrow \int_{\Omega \times \mathbb{T}^n} (a(y) (\nabla_x u(x) + \nabla_y u_1(x, y)) \cdot (\nabla_x u(x) + \nabla_y u_1(x, y))) dy dx, \end{aligned}$$

and the two-scale convergence of  $\nabla \phi^\varepsilon$  to  $\nabla_x u(x) + \nabla_y u_1(x, y)$  implies that

$$\begin{aligned} & \int_{\Omega} (a(\frac{x}{\varepsilon}) \nabla \phi^\varepsilon \cdot (\nabla_x u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}))) dx \\ & \rightarrow \int_{\Omega \times \mathbb{T}^n} (a(y) (\nabla_x u(x) + \nabla_y u_1(x, y)) \cdot (\nabla_x u(x) + \nabla_y u_1(x, y))) dx dy. \end{aligned}$$

Putting everything together gives

$$\int_{\Omega} \left| \nabla \phi_\varepsilon(x) - \left( \nabla u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}) \right) \right|^2 dx \leq C[\mathcal{L}[U, U] + \mathcal{L}[U, U] - 2\mathcal{L}[U, U]] + o(1),$$

hence

$$\int_{\Omega} \left| \nabla \phi_\varepsilon(x) - \left( \nabla u(x) + \nabla_y u_1(x, \frac{x}{\varepsilon}) \right) \right|^2 dx \rightarrow 0,$$

as  $\varepsilon \rightarrow 0$ , and we are done.

## 6 Hamilton-Jacobi equations

We will now study solutions of the initial value problem for the Hamilton-Jacobi equations of the form

$$\begin{aligned}u_t + H(\nabla u, x) &= 0 \\ u(0, x) &= u_0(x).\end{aligned}\tag{6.1}$$

In order to explain how such problems come about we need to recall some basic notions from the control theory.

### 6.1 Deterministic Optimal Control

Consider the following abstract optimization problem. Let  $y(s) : [t, T] \rightarrow \mathbb{R}^d$  denote the **state of a system** at time  $s \in [t, T]$ . This vector function could represent many things like the position and orientation of an aircraft, the amount of capital available to a government, or the wealth of an individual investor. We'll suppose that  $y(s)$  satisfies the ordinary differential equation

$$\begin{aligned}y'(s) &= f(y(s), \alpha(s)), \quad s \in [t, T] \\ y(t) &= x \in \mathbb{R}^d\end{aligned}\tag{6.2}$$

(If  $d > 1$ , this is a system of ODEs.) The function  $f(y, \alpha) : \mathbb{R}^d \times \mathbb{R}^m \rightarrow \mathbb{R}^d$  models the system dynamics. We'll suppose that  $f$  is bounded and Lipschitz continuous. The function  $\alpha(s) : [t, T] \rightarrow A$  is called a **control**. The control takes values in the set  $A$ , a compact subset of  $\mathbb{R}^m$ . The set of all possible controls or **admissible controls** will be denoted by  $\mathcal{A}_{t,T}$ :

$$\mathcal{A}_{t,T} = \{\alpha(s) : [t, T] \rightarrow A \mid \alpha(s) \text{ is measurable}\}.\tag{6.3}$$

When the dependence on  $t$  and  $T$  is clear from the context, we will simply use  $\mathcal{A}$  instead. By choosing  $\alpha$  we have some control over the course of the system  $y(t)$ . For example, in a mechanics application  $\alpha(s)$  might represent a throttle control, which determines how much thrust comes from an engine. Or, in an economics application  $\alpha$  might represent the rate at which economic resources are consumed.

So, we choose a control  $\alpha(\cdot)$  and the system evolves according to (6.2); we'd like to control the system in an optimal way, in the following sense. Suppose that the function  $g(x) : \mathbb{R}^d \rightarrow \mathbb{R}$  represents a **final payoff**; this is a reward which depends on the final state of the system at time  $T$ . Also, the function  $r(x, \alpha) : \mathbb{R}^d \times \mathbb{R}^m \rightarrow \mathbb{R}$  represents a **running payoff** or **running cost**. If  $r > 0$ , this would represent additional path-dependent payoff; if  $r < 0$ , this would represent operational costs incurred before the final payoff at time  $T$ . Given the initial state of the system  $y(t) = x$ , the optimization problem is to find an optimal control  $\alpha^*(\cdot)$  that maximizes net profit:

$$J_{x,t}(\alpha^*) = \max_{\alpha(\cdot) \in \mathcal{A}} J_{x,t}(\alpha) = \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^T r(y(s), \alpha(s)) ds + g(y(T)) \right]\tag{6.4}$$

If  $r < 0$ , this may be thought of as the optimal balance between payoff and running costs. Although, we may be able to control the system (by choosing  $\alpha$ ) so that the final payoff

$g(y(T))$  is large, it may be very expensive to arrive at this state. So, we want to find the optimal control that balances these competing factors.

Even if an optimal control does not exist, we may study the function

$$u(x, t) = \max_{\alpha(\cdot) \in \mathcal{A}} J_{x,t}(\alpha) = \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^T r(y(s), \alpha(s)) ds + g(y(T)) \right] \quad (6.5)$$

This function is called the **value function** associated with the control problem. It depends on  $x$  and  $t$  through the initial conditions defining  $y(s)$ . There are many interesting mathematical questions related to this optimization problem. For example:

1. For given  $(x, t)$ , is there an optimal control  $\alpha^*$ ?
2. If so, how can we compute it?
3. How does the value function  $u$  depend on  $x$  and  $t$ ? Does it satisfy a PDE?

## Some Examples

Here are two examples which fit into this abstract framework.

**Example 1:** Suppose you want to drive a boat from position  $x_0$  at time  $t$  to a position  $x_f$  at time  $T$ . Let  $x(s)$  denote the position of the boat,  $v(s)$  denote the velocity, then a simple model for the boat dynamics might be

$$\begin{aligned} x'(s) &= v(s) \\ v'(s) &= \frac{\alpha(s)}{(m_1 + m(s))} - \beta(v(s)) \\ m'(s) &= -k|\alpha(s)| \end{aligned}$$

Here  $m(s)$  is the mass of the boat's fuel,  $m_1$  is the boat's dry weight, the function  $\beta(v) \geq 0$  models drag as the boat moves through the water. The vector  $\alpha(s)$  represents a throttle and direction control, and its magnitude is proportional to the rate of fuel consumption ( $k$  is a proportionality constant). The acceleration is also proportional to the throttle control parameter.

How should we steer the boat in order to minimize fuel consumption? In this setting, the system state  $y(s)$  is the vector  $y(s) = (x(s), v(s), m(s))$ . One way to model this problem would be to find

$$\max_{\alpha} J_{x_0,t}(\alpha) = \max_{\alpha} [m(T) + p(x(T))] \quad (6.6)$$

Here  $p \leq 0$  might be a function satisfying  $p(x) = 0$  if  $x$  is close to  $x_f$  and  $p(x) \ll -1$  if  $x$  is far from  $x_f$ . So, although we don't need to land precisely at  $x_f$ , there is a big penalty for leaving the boat far from  $x_f$ . There is no "running cost" in this model. Notice that it is possible for  $m(s)$  to become negative, which is non-physical. We could fix this modeling issue by modifying the equations appropriately or by applying an additional **state constraint** of the form  $0 \leq m(s)$ .

**Example 2:** Here is a variant of a classic example studied by F. P. Ramsey (see *The Economic Journal*, Vol. 38, No. 152, 1928). The problem is to determine how much of a nation's resources should be saved and how much should be consumed. Let  $c(s)$  denote the rate of capital consumption, let  $p(s)$  denote the rate of capital production, and let  $k(s)$  denote the amount of capital at time  $s$ . Then the rate change in capital is the difference between the rates of production and consumption:

$$k'(s) = p(s) - c(s). \quad (6.7)$$

Suppose that the production is related to capital and consumption as  $p(s) = P(c(s), k(s))$ . Suppose also that consumption is related to capital according to  $c(s) = \alpha(s)C(k(s))$ , where  $\alpha(s)$  is a control. Therefore,

$$k'(s) = P(\alpha(s)C(k(s)), k(s)) - \alpha(s)C(k(s)). \quad (6.8)$$

Given current level of capital  $k(t) = k_0$ , we'd like to choose a level of consumption (by choosing  $\alpha$ ) which maximizes the total utility; this goal might be modeled by the optimal control problem

$$\max_{\alpha} J_{k_0, t}(\alpha) = \max_{\alpha} \left[ \int_t^T U(c(s)) ds + U_T(k(T)) \right] = \max_{\alpha} \left[ \int_t^T U(\alpha(s)C(k(s))) ds + U_T(k(T)) \right].$$

Here  $U$  is a utility function, and  $U_T$  models some payoff representing the utility of having left-over capital  $k(T)$  at time  $T$ .

## 6.2 The Dynamic Programming Principle

**Theorem 6.1** *Let  $u(x, t)$  be the value function defined by (6.5). If  $t < \tau \leq T$ , then*

$$u(x, t) = \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^{\tau} r(y(s), \alpha(s)) ds + u(y(\tau), \tau) \right]. \quad (6.9)$$

The relation (6.44) is called the **Dynamic Programming Principle**, and it is a fundamental tool in the analysis of optimal control problems. It says that if we know the value function at time  $\tau > t$ , we may determine the value function at time  $t$  by optimizing from time  $t$  to time  $\tau$  and using  $u(\cdot, \tau)$  as the payoff. Roughly speaking, this is reminiscent of the Markov property of a stochastic process, in the sense that if we know  $u(x, \tau)$  we can determine  $u(\cdot, t)$  for  $t < \tau$  without any other information about the control problem beyond time  $\tau$  (ie. times  $s \in [\tau, T]$ ). (More precisely, it means that  $u(x, t)$  satisfies what is called a semi-group property.)

**Proof of Theorem 6.1:** At the heart of this proof of the Dynamic Programming Principle is the observation that any admissible control  $\alpha \in \mathcal{A}_{t, T}$  is the combination of a control in  $\mathcal{A}_{t, \tau}$  with a control in  $\mathcal{A}_{\tau, T}$ . We will express this relationship as

$$\mathcal{A}_{t, T} = \mathcal{A}_{t, \tau} \oplus \mathcal{A}_{\tau, T}. \quad (6.10)$$

This notation  $\oplus$  means that if  $\alpha_t(s) \in \mathcal{A}_{t, \tau}$  and  $\alpha_{\tau}(s) \in \mathcal{A}_{\tau, T}$ , then the control defined by splicing  $\alpha_t$  and  $\alpha_{\tau}$  according to

$$\alpha(s) = (\alpha_t \oplus \alpha_{\tau})(s) := \begin{cases} \alpha_t(s), & s \in [t, \tau] \\ \alpha_{\tau}(s), & s \in [\tau, T] \end{cases} \quad (6.11)$$

is an admissible control in  $\mathcal{A}_{t,T}$ . On the other hand, if we have  $\alpha \in \mathcal{A}_{t,T}$ , then by restricting the domain of  $\alpha$  to  $[t, \tau]$  we obtain an admissible control in  $\mathcal{A}_{t,\tau}$ . Similarly, by restricting the domain of  $\alpha$  to  $[\tau, T]$  we obtain an admissible control in  $\mathcal{A}_{\tau,T}$ .

The function  $u$  is defined as

$$\begin{aligned} u(x, t) &= \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^T r(y(s), \alpha(s)) ds + g(y(T)) \right] \\ &= \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^\tau r(y(s), \alpha(s)) ds + \int_\tau^T r(y(s), \alpha(s)) ds + g(y(T)) \right]. \end{aligned}$$

Notice that the first integral on the right depends only on  $y$  and  $\alpha$  up to time  $\tau$ , while the last two terms depend on the values of  $y$  and  $\alpha$  after time  $\tau$ . Since a control  $\alpha \in \mathcal{A}_{t,T}$  may be decomposed as  $\alpha = \alpha_1 \oplus \alpha_2$  with  $\alpha_1 \in \mathcal{A}_{t,\tau}$  and  $\alpha_2 \in \mathcal{A}_{\tau,T}$ , we may maximize over each component in the decomposition:

$$\begin{aligned} u(x, t) &= \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^\tau r(y(s), \alpha(s)) ds + \int_\tau^T r(y(s), \alpha(s)) ds + g(y(T)) \right] \\ &= \max_{\alpha_1 \in \mathcal{A}_{t,\tau}, \alpha_2 \in \mathcal{A}_{\tau,T}, \alpha = \alpha_1 \oplus \alpha_2} \left[ \int_t^\tau r(y(s), \alpha(s)) ds + \int_\tau^T r(y(s), \alpha(s)) ds + g(y(T)) \right]. \end{aligned}$$

On the right hand side, the system state  $y(t)$  is determined by (6.2) with  $\alpha = \alpha_1 \oplus \alpha_2 \in \mathcal{A}_{t,T}$ . Therefore, we may decompose the system state as  $y(s) = y_1 \oplus y_2$  where  $y_1(s) : [t, \tau] \rightarrow \mathbb{R}^d$  and  $y_2(s) : [\tau, T] \rightarrow \mathbb{R}^d$  are defined by

$$\begin{aligned} y_1'(s) &= f(y_1(s), \alpha_1(s)), \quad s \in [t, \tau] \\ y_1(t) &= x \end{aligned}$$

and

$$\begin{aligned} y_2'(s) &= f(y_2(s), \alpha_2(s)), \quad s \in [\tau, T] \\ y_2(\tau) &= y_1(\tau) = y(\tau). \end{aligned}$$

Here we use  $\oplus$  to denote the splicing or gluing of  $y_1$  and  $y_2$  to create  $y(t) : [t, T] \rightarrow \mathbb{R}^d$ . Therefore,

$$u(x, t) = \max_{\alpha_1 \in \mathcal{A}_{t,\tau}} \max_{\alpha_2 \in \mathcal{A}_{\tau,T}, y_2(\tau) = y_1(\tau)} \left[ \int_t^\tau r(y_1(s), \alpha_1(s)) ds + \int_\tau^T r(y_2(s), \alpha_2(s)) ds + g(y_2(T)) \right],$$

where the initial point for  $y_2(\tau)$  is  $y_2(\tau) = y_1(\tau)$ . Observe that  $y_1$  depends only on  $x$  and  $\alpha_1$ , not on  $y_2$  or  $\alpha_2$ . Since the first integral depends only on  $\alpha_1$  and  $y_1$ , this may be rearranged as

$$\begin{aligned} u(x, t) &= \max_{\alpha_1 \in \mathcal{A}_{t,\tau}} \max_{\alpha_2 \in \mathcal{A}_{\tau,T}, y_2(\tau) = y_1(\tau)} \left[ \int_t^\tau r(y_1(s), \alpha_1(s)) ds + \int_\tau^T r(y_2(s), \alpha_2(s)) ds + g(y_2(T)) \right] \\ &= \max_{\alpha_1 \in \mathcal{A}_{t,\tau}} \left[ \int_t^\tau r(y_1(s), \alpha_1(s)) ds + \max_{\alpha_2 \in \mathcal{A}_{\tau,T}, y_2(\tau) = y_1(\tau)} \left( \int_\tau^T r(y_2(s), \alpha_2(s)) ds + g(y_2(T)) \right) \right] \\ &= \max_{\alpha_1 \in \mathcal{A}_{t,\tau}} \left[ \int_t^\tau r(y_1(s), \alpha_1(s)) ds + u(y_1(\tau), \tau) \right] \quad (\text{using the definition of } u) \\ &= \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^\tau r(y(s), \alpha(s)) ds + u(y(\tau), \tau) \right] \end{aligned} \tag{6.12}$$

This completes the proof.  $\square$

Notice that in this proof we have not assumed that an optimal control exists.

### 6.3 The Hamilton-Jacobi-Bellman Equation

How does the value function depend on  $x$  and  $t$ ? Is it continuous in  $(x, t)$ ? Is it differentiable? Does it satisfy a PDE? Unfortunately, the value function may not be differentiable, even in simple examples! Here is one interesting example of this fact. Suppose that  $f(x, \alpha) = \alpha$ ,  $g \equiv 0$ , and  $r(x, \alpha)$  is defined by

$$r(x, \alpha) = -\mathbb{I}_D(x) = \begin{cases} -1, & x \in D \\ 0, & x \in \mathbb{R}^d \setminus D \end{cases} \quad (6.13)$$

where  $D \subset \mathbb{R}^d$  is some bounded set. Suppose that the set of admissible controls is defined by (6.49) with  $A = \{|\alpha| \leq 1\}$ . In this case,  $y'(s) = \alpha(s)$  and  $|y'(s)| \leq 1$ . Therefore, the value function may be written as

$$u(x, t) = \max_{y: [t, T] \rightarrow \mathbb{R}^d, |y'| \leq 1, y(t) = x} \left[ \int_t^T -\mathbb{I}_D(y(s)) ds \right]. \quad (6.14)$$

Clearly  $u(x, t) \leq 0$ , and the optimum is obtained by paths that spend the least amount of time in the set  $D$ . If  $x \in \mathbb{R}^d \setminus D$ , then  $u(x, t) = 0$ , because we could take  $y(s) = x$  for all  $s \in [t, T]$ . In this case, the system state doesn't change, so the integral is zero, which is clearly optimal. On the other hand, if  $x \in D$  then the optimal control moves  $y(s)$  to  $\mathbb{R}^d \setminus D$  as quickly as possible and then stays outside  $D$ . Since  $|y'(s)| \leq 1$ , this implies that the value function is given explicitly by

$$u(x, t) = -\min((T - t), \text{dist}(x, \mathbb{R} \setminus D)) \quad (6.15)$$

where

$$\text{dist}(x, \mathbb{R} \setminus D) = \inf_{y \in \mathbb{R} \setminus D} |x - y|, \quad (6.16)$$

is the Euclidean distance from  $x$  to the outside of  $D$ . Albeit continuous, this function may not be differentiable! For example, suppose that  $D = \{(x_1, x_2) \in \mathbb{R}^2 \mid x_1^2 + x_2^2 \leq 1\}$  is the unit disk. In this case,

$$u(x, t) = \begin{cases} |x| - 1, & |x| \leq 1 \\ 0, & |x| \geq 1 \end{cases} \quad (6.17)$$

for  $t \leq T - 1$ . Thus  $u(x, t)$  is not differentiable at the origin  $x = (x_1, x_2) = (0, 0)$  for  $t < T - 1$ .

So, in general, the value function may not be differentiable. However, one can still derive a PDE satisfied by the value function. If the value function is differentiable, this equation is satisfied in the classical sense. At points where the value function is not differentiable, one can show that the value function (assuming it is at least continuous) satisfies the PDE in a weaker sense. This weaker notion of "solution" is called a "viscosity solution" of the PDE. For the moment, we will formally compute as if the value function were actually differentiable. We will come back to the weak solutions later.

For now, let us use the Dynamic Programming Principle to formally derive an equation solved by the value function  $u(x, t)$ . The Dynamic Programming Principle does not require differentiability of the value function; however, in our computations we assume that the value function is continuous and differentiable with respect to both  $x$  and  $t$ . The Dynamic Programming Principle tells us that

$$u(x, t) = \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^\tau r(y(s), \alpha(s)) ds + u(y(\tau), \tau) \right]. \quad (6.18)$$

To formally derive a PDE for  $u$ , we let  $h \in (0, T - t)$  and set  $\tau = t + h < T$ , then

$$u(x, t) = \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^{t+h} r(y(s), \alpha(s)) ds + u(y(t+h), t+h) \right]. \quad (6.19)$$

We'll assume that nearly optimal controls are approximately constant for  $s \in [t, t+h]$ .

First, consider the term  $u(y(t+h), t+h)$ . From the chain rule and our assumption that  $u$  is continuously differentiable in  $x$  and  $t$ , we conclude that

$$u(y(t+h), t+h) = u(y(t), t) + hy'(t) \cdot \nabla u(y(t), t) + hu_t(y(t), t) + o(h). \quad (6.20)$$

Recall that a function  $\phi(h)$  is said to be  $o(h)$  ("little oh of  $h$ ") if  $\lim_{h \rightarrow 0} (\phi(h)/h) = 0$ . So, (6.20) says that  $u(y(t+h), t+h)$  is equal to a linear function of  $h$  plus something that is  $o(h)$  (i.e. higher order than  $h$ , but not necessarily  $O(h^2)$ ). Therefore,

$$\begin{aligned} u(y(t+h), t+h) &= u(y(t), t) + hf(y(t), \alpha(t)) \cdot \nabla u(y(t), t) + hu_t(y(t), t) + o(h) \\ &= u(x, t) + hf(x, \alpha(t)) \cdot \nabla u(x, t) + hu_t(x, t) + o(h). \end{aligned} \quad (6.21)$$

Now, plug this into (6.19):

$$u(x, t) = \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^{t+h} r(y(s), \alpha(s)) ds + u(x, t) + hf(x, \alpha(t)) \cdot \nabla u(x, t) + hu_t(x, t) + o(h) \right]. \quad (6.22)$$

The term  $u(x, t)$  may be pulled out of the maximum, so that it cancels with the left hand side:

$$0 = hu_t(x, t) + o(h) + \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^{t+h} r(y(s), \alpha(s)) ds + hf(x, \alpha(t)) \cdot \nabla u(x, t) \right]. \quad (6.23)$$

Now, divide by  $h$  and let  $h \rightarrow 0$ .

$$0 = u_t(x, t) + \frac{o(h)}{h} + \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \frac{1}{h} \int_t^{t+h} r(y(s), \alpha(s)) ds + f(x, \alpha(t)) \cdot \nabla u(x, t) \right]. \quad (6.24)$$

If  $\alpha(s)$  is continuous at  $t$ , then as  $h \rightarrow 0$ ,

$$\lim_{h \rightarrow 0} \frac{1}{h} \int_t^{t+h} r(y(s), \alpha(s)) ds = r(y(t), \alpha(t)) = r(x, \alpha(t)) \quad (6.25)$$

So, if the nearly optimal controls are approximately constant for  $s \in [t, t + h]$ , then by letting  $h \rightarrow 0$  in (6.24) we conclude that

$$u_t(x, t) + \max_{a \in A} [r(x, a) + f(x, a) \cdot \nabla u(x, t)] = 0, \quad x \in \mathbb{R}^d, t < T. \quad (6.26)$$

This equation is called the **Hamilton-Jacobi-Bellman equation**. The function  $u(x, t)$  also satisfies the terminal condition

$$u(x, T) = g(x). \quad (6.27)$$

Notice that the HJB equation is a first-order, fully nonlinear equation, having the form

$$u_t + H(\nabla u, x) = 0$$

where the function  $H$  is defined by

$$H(p, x) = \max_{a \in A} [r(x, a) + f(x, a) \cdot p], \quad p \in \mathbb{R}^d, \quad (6.28)$$

and is sometimes called the **Hamiltonian**.

In addition to telling us how the value function depends on  $x$  and  $t$ , this PDE suggests what the optimal control should be. Suppose  $u(x, t)$  is differentiable and solves the PDE in the classical sense. Then

$$u_t + H(\nabla u, x) = 0, \quad (6.29)$$

where  $H(p, x)$  is defined by (6.28). Then the optimal control is computed by finding  $(y^*(s), \alpha^*(s))$  which satisfies

$$H(\nabla u(y^*(s), s), y^*(s)) = r(y^*(s), \alpha^*(s)) + f(y^*(s), \alpha^*(s)) \cdot \nabla u(y^*(s), s), \quad (6.30)$$

and

$$\begin{aligned} \frac{d}{dt} y^*(s) &= f(y^*(s), \alpha^*(s)), \quad s > t \\ y^*(t) &= x. \end{aligned} \quad (6.31)$$

## Examples

**Example 1:** In this example we want to maximize utility from consumption of our resources over time interval  $[t, T]$ . If  $c(s)$  is the rate of consumption, then we model the utility gained from this consumption as

$$U = \int_t^T e^{-\mu s} \psi(c(s)) ds. \quad (6.32)$$

For example, we might choose  $\psi$  to be an increasing, concave function of  $c$  like  $\psi(c) = c^\nu$  for some power  $\nu \in (0, 1)$ . The factor  $e^{-\mu s}$  is a discount factor. Let us suppose that the rate of consumption is  $c(s) = \alpha(s)y(s)$ , where  $y$  is our total wealth and  $\alpha$  is a control. So,  $\alpha$  is approximately the proportion of total wealth consumed in a unit of time. Instead of consuming resources, we might invest them in such a way that our total wealth satisfies the ode

$$y'(s) = qy(s) - c(s) = qy(s) - \alpha(s)y(s) \quad (6.33)$$

Without consumption, our wealth would grow exponentially at rate  $q > 0$ . So, investing some of our wealth might allow us to consume more in the long run. The control  $\alpha(s)$  should satisfy some realistic constraints, for example  $\alpha \in A = [0, \bar{a}]$ . Thus, if  $x$  denotes the current wealth  $y(t) = x$ , the control problem is to find

$$u(x, t) = \max_{\alpha \in A} \left[ \int_t^T e^{-\mu s} (\alpha(s)y(s))^\nu ds + g(y(T)) \right] \quad (6.34)$$

The function  $g$  models some utility of leaving an amount of unused wealth  $y(T)$  at the final time. In this example, we can actually compute  $y(s)$  explicitly:

$$y(s) = y(t)e^{q(t-s) - \int_t^s \alpha(\tau) d\tau} \quad (6.35)$$

To determine the associated HJB equation, notice that  $f(x, a) = (q - a)x$  and  $r(x, a, s) = e^{-\mu s}(ax)^\nu$ . Therefore, we expect the value function to satisfy

$$u(x, t) + \max_{a \in A} [e^{-\mu s}(ax)^\nu + (q - a)x \cdot \nabla u] = 0, \quad (6.36)$$

perhaps in a weak sense (viscosity solutions) rather than a classical sense.

**Example 2:** Here's an example from engineering. Suppose that  $r \equiv 0$ ,  $f(x, a) = -v(x) + a$ , and  $A = \{|a| \leq \mu_0\}$ . In this case, the HJB equation is

$$u_t(x, t) + \max_{|a| \leq \mu_0} [-v(x) \cdot \nabla u(x, t) + a \cdot \nabla u(x, t)] = 0, \quad x \in \mathbb{R}^d, t < T. \quad (6.37)$$

It is easy to see that the optimal  $a$  is  $a = \mu_0(\nabla u)/|\nabla u|$ , so that the equation becomes

$$u_t - v(x) \cdot \nabla u + \mu_0 |\nabla u| = 0 \quad (6.38)$$

The function  $G(x, t) = u(x, T - t)$  satisfies

$$G_t + v(x) \cdot \nabla G = \mu_0 |\nabla G|. \quad (6.39)$$

In the combustion community, this equation is called the ‘‘G-equation’’ and is used in computational models of premixed turbulent combustion. The level set  $\{G = 0\}$  is considered to be the flame surface. The parameter  $\mu_0$  corresponds to the laminar flame speed (i.e. the flame speed without the turbulent velocity field  $v$ ).

## Infinite Time Horizon

So far we have considered a deterministic control problem with **finite time horizon**. This means that the optimization involves a finite time interval and may involve a terminal payoff. One might also consider an optimization posed on an infinite time interval. Suppose that  $y : [t, \infty) \rightarrow \mathbb{R}^d$  satisfies

$$\begin{aligned} y'(s) &= f(y(s), \alpha(s)), \quad s \in [t, \infty) \\ y(t) &= x \in \mathbb{R}^d. \end{aligned} \quad (6.40)$$

Now the domain for the control is also  $[t, \infty)$ . We'll use  $\mathcal{A} = \mathcal{A}_{t, \infty}$  for the set of admissible controls. For  $x \in \mathbb{R}^d$ , define the value function

$$u(x, t) = \max_{\alpha(\cdot) \in \mathcal{A}} J_{x, t}(\alpha) = \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_t^\infty e^{-\lambda s} r(y(s), \alpha(s)) ds \right]. \quad (6.41)$$

The exponential term in the integral is a discount factor; without it, the integral might be infinite. Notice that there is no terminal payoff, only running payoff. This optimal control problem is said to involve an **infinite time horizon**. Notice also that the value functions depends on  $t$  in a trivial way. Since  $r$  and  $f$  do not depend on  $t$ , we may change variables to see that

$$u(x, t) = e^{-\lambda t} u(x, 0). \quad (6.42)$$

So, to find  $u(x, t)$  it suffices to compute

$$u(x) = \max_{\alpha(\cdot) \in \mathcal{A}} J_x(\alpha) = \max_{\alpha(\cdot) \in \mathcal{A}} \left[ \int_0^\infty e^{-\lambda s} r(y(s), \alpha(s)) ds \right] \quad (6.43)$$

where  $\mathcal{A} = \mathcal{A}_{0, \infty}$ .

**Theorem 6.2 (Dynamic Programming Principle)** *Let  $u(x)$  be the value function defined by (6.43). For any  $x \in \mathbb{R}^d$  and  $h > 0$ ,*

$$u(x) = \max_{\alpha(\cdot) \in \mathcal{A}_{0, h}} \left[ \int_0^h e^{-\lambda s} r(y(s), \alpha(s)) ds + e^{-\lambda h} u(y(h)) \right] \quad (6.44)$$

**Proof:** Exercise.

Using the Dynamic Programming Principle, one can formally derive the HJB equation for the infinite horizon control problem. The equation is:

$$-\lambda u + \max_{a \in A} [r(x, a) + f(x, a) \cdot \nabla u] = 0 \quad (6.45)$$

which has the form

$$-\lambda u + H(\nabla u, x) = 0 \quad (6.46)$$

with the Hamiltonian  $H(p, x)$  defined by

$$H(p, x) = \max_{a \in A} [r(x, a) + f(x, a) \cdot p] \quad (6.47)$$

**Exercise:** Check these computations.

## 6.4 Brief Introduction to Stochastic Optimal Control

Thus far, we have considered deterministic optimal control in which the dynamic behaviour of the system state is deterministic. In a stochastic optimal control problem, the system state  $y(s)$  is a stochastic process. Consequently, the controls also will be stochastic, since we may want to steer the system in a manner that depends on the system's stochastic trajectory. To this end, we now suppose that the **system state**  $Y_s(\omega) : [t, T] \times \Omega \rightarrow \mathbb{R}^d$  now satisfies the stochastic differential equation (or system of equations)

$$\begin{aligned} dY_s &= f(Y_s, \alpha_s, s)ds + \sigma(Y_s, \alpha_s, s)dB_s, \quad s \geq t \\ Y_t &= x, \quad a.s. \end{aligned} \tag{6.48}$$

where  $B_s$  is a  $n$ -dimensional Brownian motion defined on probability space  $(\Omega, \mathcal{F}, \{\mathcal{F}_s\}_{s \geq t}, P)$ , and  $\sigma$  is a  $d \times n$  matrix.

We control the system state through the **control process**  $\alpha_s(\omega) : [t, T] \times \Omega \rightarrow \mathbb{R}^m$  which is adapted to the filtration  $\{\mathcal{F}_s\}_{s \geq t}$ . The set of admissible controls is now

$$\mathcal{A}_{t,T} = \{\alpha_s(\omega) : [t, T] \times \Omega \rightarrow A \mid \alpha_s \text{ is adapted to the filtration } \{\mathcal{F}_s\}_{s \geq t}\}. \tag{6.49}$$

The assumption that the controls are adapted means that we cannot look into the future; the control can only be chosen on the basis of information known up to the present time. Supposing that  $\sigma$  and  $f$  satisfy the usual bounds and continuity conditions, the stochastic process  $Y_s(\omega)$  is uniquely determined by the initial condition  $Y_t = x$  and the control process  $\alpha_s(\omega)$ .

Given a time  $T > t$ , the abstract stochastic optimal control problem is to maximize

$$\max_{\alpha \in \mathcal{A}_{t,T}} J_{x,t}(\alpha(\cdot)) = \max_{\alpha \in \mathcal{A}_{t,T}} E \left[ \int_t^T r(Y_s, \alpha_s, s) ds + g(Y_T) \mid Y_t = x \right] \tag{6.50}$$

As before, the function  $r(y, \alpha, s)$  represents a running payoff (or running cost, if  $r < 0$ ), and  $g$  represents a terminal payoff (or terminal cost, if  $g < 0$ ). Since the system state is a stochastic process, the net payoff is a random variable, and our goal is to maximize the expected payoff. Even if an optimal control process does not exist, we may define the value function to be

$$u(x, t) = \max_{\alpha \in \mathcal{A}_{t,T}} E \left[ \int_t^T r(Y_s, \alpha_s, s) ds + g(Y_T) \mid Y_t = x \right] \tag{6.51}$$

Notice that the value function is *not* random.

### Variations

There are variations on this theme. For example, we might add the possibility of a payoff based on a stopping criteria. In this case, we want to maximize:

$$\max_{\alpha \in \mathcal{A}_{t,T}} E \left[ \int_t^{\gamma \wedge T} r(Y_s, \alpha_s, s) ds + g(Y_T)\mathbb{I}_{\{\gamma \geq T\}} + h(Y_\gamma)\mathbb{I}_{\{\gamma < T\}} \mid Y_t = x \right] \tag{6.52}$$

Here, the random variable  $\gamma(\omega)$  is a stopping time. The function  $h$  represents some payoff that is incurred if  $\gamma < T$  (or, this may represent a penalty if  $h < 0$ ).

In (6.51) the time horizon is finite. One could also pose an optimal control problem on an infinite time horizon. For example, one might consider maximizing

$$\max_{\alpha \in \mathcal{A}} J_{x,t}(\alpha(\cdot)) = \max_{\alpha \in \mathcal{A}} E \left[ \int_t^\gamma e^{-\lambda s} r(Y_s, \alpha_s) ds + e^{-\lambda \gamma} h(Y_\gamma) \right] \quad (6.53)$$

where  $\gamma$  is a stopping time.

## 6.5 Dynamic Programming Principle for Stochastic Control

For the stochastic control problem there is a Dynamic Programming Principle that is analogous to the DPP for deterministic control. Using the Markov Property of the stochastic process  $Y_t$ , one can easily prove the following:

**Theorem 6.3** *Let  $u(x, t)$  be the value function defined by (6.51). If  $t < \tau \leq T$ , then*

$$u(x, t) = \max_{\alpha \in \mathcal{A}_{t,\tau}} E \left[ \int_t^\tau r(Y_s, \alpha_s, s) ds + u(Y_\tau, \tau) \mid Y_t = x \right] \quad (6.54)$$

**Proof:** Exercise. The idea is the same as in the case of deterministic control. Split the integral into two pieces, one over  $[t, \tau]$  and the other over  $[\tau, T]$ . Then condition on  $\mathcal{F}_\tau$  and use the Markov property, so that the second integral and the payoff may be expressed in terms of  $u(Y_\tau, \tau)$ .  $\square$

## 6.6 HJB equation

Using the Dynamic Programming Principle, one can formally derive a PDE for the value function  $u(x, t)$ . As in the case of deterministic optimal control, one must assume that the value function is sufficiently smooth. Because the dynamics are stochastic, we want to apply Itô's formula in the way that we used the chain rule to derive the HJB equation for deterministic control. Thus, this formal computation requires that the value function be twice differentiable.

From Itô's formula we see that

$$\begin{aligned} u(Y_\tau, \tau) - u(x, t) &= \int_t^\tau [u_t(Y_s, s) + f(Y_s, \alpha_s, s) \cdot \nabla u(Y_s, s)] ds \\ &\quad + \int_t^\tau \frac{1}{2} \sum_k \sum_{i,j} u_{x_i x_j}(Y_s, s) \sigma^{jk}(Y_s, \alpha_s, s) \sigma^{ik}(Y_s, \alpha_s, s) ds \\ &\quad + \int_t^\tau (\nabla u(Y_s, s))^T \sigma(Y_s, \alpha_s, s) dB_s \\ &= \int_t^\tau u_t(Y_s, s) + \mathcal{L}^\alpha u(Y_s, s) ds + \int_t^\tau (\nabla u(Y_s, s))^T \sigma(Y_s, \alpha_s, s) dB_s \end{aligned} \quad (6.55)$$

where  $\mathcal{L}$  is the second order differential operator

$$\begin{aligned} \mathcal{L}^\alpha u &= f(y, \alpha, s) \cdot \nabla u(y, s) + \frac{1}{2} \sum_k \sum_{i,j} u_{y_i y_j}(y, s) \sigma^{jk}(y, \alpha_s, s) \sigma^{ik}(y, \alpha_s, s) \\ &= f(y, \alpha, s) \cdot \nabla u(y, s) + \frac{1}{2} \text{tr}(D^2 u(y, s) \sigma(y, \alpha, s) \sigma^T(y, \alpha, s)) \end{aligned} \quad (6.56)$$

and  $D^2u$  is the matrix of second partial derivatives. Now we plug this into the DPP relation (6.54) and use the fact the martingale term in (6.55) has zero mean. We obtain:

$$\begin{aligned} 0 &= \max_{\alpha \in \mathcal{A}_{t,\tau}} E \left[ \int_t^\tau r(Y_s, \alpha_s, s) ds + u(Y_\tau, \tau) - u(x, t) \mid Y_t = x \right] \\ &= \max_{\alpha \in \mathcal{A}_{t,\tau}} E \left[ \int_t^\tau r(Y_s, \alpha_s, s) ds + \int_t^\tau u_t(Y_s, s) + \mathcal{L}^\alpha u(Y_s, s) ds \mid Y_t = x \right] \end{aligned} \quad (6.57)$$

Finally, let  $\tau = t + h$ , divide by  $h$  and let  $h \rightarrow 0$ , as in the deterministic case. We formally obtain the HJB equation

$$u_t(x, t) + \max_{a \in A} [r(x, a, t) + \mathcal{L}^a u(x, t)] = 0. \quad (6.58)$$

This may be written as

$$u_t(x, t) + \max_{a \in A} \left[ r(x, a, t) + \frac{1}{2} \text{tr}(D^2u(x, t) \sigma(x, a, t) \sigma^T(x, a, t)) + f(x, a, t) \cdot \nabla u(x, t) \right] = 0 \quad (6.59)$$

which is, in general, a fully-nonlinear, *second order* equation of the form

$$u_t + H(D^2u, Du, x, t) = 0 \quad (6.60)$$

Notice that the equation is deterministic. The set of possible control values  $A \subset \mathbb{R}^m$  is a subset of Euclidean space, and the maximum in the HJB equation (6.68) is over this deterministic set, not over the set  $\mathcal{A}$ .

## HJB for the infinite horizon problem

Deriving the HJB for the infinite horizon problem is similar. Suppose  $r = r(y, a)$  and  $f = f(y, a)$ . Suppose that  $u$  is the value function defined by

$$u(x) = \max_{\alpha \in \mathcal{A}} E \left[ \int_0^\infty e^{-\lambda s} r(Y_s, \alpha_s) ds \mid Y_0 = x \right] \quad (6.61)$$

Then the Dynamic Programming Principle shows that for any  $\tau > 0$

$$u(x) = \max_{\alpha \in \mathcal{A}} E \left[ \int_0^\tau e^{-\lambda s} r(Y_s, \alpha_s) ds + e^{-\lambda \tau} u(Y_\tau) \mid Y_0 = x \right] \quad (6.62)$$

Using Itô's formula as before, we formally derive the second order equation

$$-\lambda u(x) + \max_{a \in A} [r(x, a) + \mathcal{L}^a u(x)] = 0 \quad (6.63)$$

## 6.7 Examples

**Example 1:** In this example, we consider the problem of portfolio optimization. We already considered the problem of maximizing utility from consumption of resources. In that model,

we assumed that in the absence of consumption the individual's total wealth grows exponentially according to a deterministic rate (interest rate). Now we suppose that the individual may invest money in various asset classes and or consume resources. For simplicity we suppose that the individual may invest in either a stock (risky) or a bond (risk free growth). Suppose that the stock and bond satisfy the equations

$$\begin{aligned} db_s &= rb_s ds \\ dP_s &= \mu P_s ds + \sigma P_s dB_s \end{aligned} \quad (6.64)$$

Here  $B_s$  is a standard Brownian motion. The individual may decide how much money to consume, how much money to invest in stock, and how much money to invest in bonds. Let  $c_s$  denote the consumption rate. Let  $\pi_s$  denote the proportion of investments to put in stocks. Now, if  $Y_s$  is the individual's wealth at time  $s$ ,  $Y_s$  satisfies the equation

$$\begin{aligned} dY_s &= (1 - \pi_s)rY_s ds + \pi_s\mu Y_s ds + \pi_s\sigma Y_s dB_s - c_s ds \\ &= [((1 - \pi_s)r + \pi_s\mu)Y_s - c_s] ds + \pi_s\sigma Y_s dB_s \end{aligned} \quad (6.65)$$

Lets suppose that  $c_s = \eta_s Y_s$  with  $\eta \geq 0$ , so that the consumption rate cannot be negative. Then the control has two components:  $\alpha_s = (\eta_s, \pi_s)$ . Now our goal is to maximize the discounted utility

$$J_x(\alpha(\cdot)) := \int_0^\infty e^{-\lambda s} U(c_s) ds = \int_0^\infty e^{-\lambda s} U(\eta_s Y_s) ds \quad (6.66)$$

where  $U(c)$  is a utility function, typically increasing and concave. So the value function is

$$u(x) = \max_{\alpha \in \mathcal{A}} \int_0^\infty e^{-\lambda s} U(\eta_s Y_s) ds \quad (6.67)$$

Notice that  $u$  depends on  $\pi_s$  through the definition of  $Y_s$ .

The HJB equation has the form

$$-\lambda u(x) + \max_{\eta, \pi} \left[ U(\eta y) + ((1 - \pi)r + \pi\mu - \eta)xu_x(x) + \frac{\pi^2 \sigma^2}{2} x^2 u_{xx}(x) \right] = 0. \quad (6.68)$$

This example is based on the paper by Robert Merton, *Optimal Consumption and Portfolio Rules in a Continuous-Time Model*, J. Econ. Theory, **3**, 1971 pp. 373-413. The model can be solved explicitly. See lecture notes by M. Soner for details.

**Example 2:** This example comes from the interesting paper by Leung, Sircar, and Zariphopoulou cited above. An investor want to optimize her portfolio, investing in either stock or bond, with no consumption. The problem is that the company in which she invests may default. If this occurs, she must sell her stock at the market price and put her money in the bond (she can't re-invest in another stock for the time interval under consideration). The stock price is modeled as:

$$dP_s = \mu P_s ds + \sigma P_s dB_s^1 \quad (6.69)$$

where  $B_s^1$  is a Brownian motion. The value of the firms assets is modeled as

$$dV_s = \nu V_s ds + \eta V_s (\rho dB_s^1 + \rho' dB_s^2) \quad (6.70)$$

where  $B_s^2$  is an independent Brownian motion, and  $\rho \in (-1, 1)$ ,  $\rho' = \sqrt{1 - \rho^2}$ . Default occurs if the firm's price  $P_s$  drops below a boundary  $\tilde{D}_s = D e^{-\beta(T-s)}$ ,  $s \in [0, T]$ ,  $D > 0$ . Choosing to invest a ratio  $\pi_s$  in the stock (as in the preceding example), the investor's wealth is  $Y_s$  where

$$dY_s = (1 - \pi_s)rY_s ds + \pi_s\mu Y_s ds + \pi_s\sigma Y_s dB_s \quad (6.71)$$

and  $r$  is the interest rate for the bond. The initial condition is  $Y_t = y$ . We define the default time, a stopping time, by

$$\gamma_t = \inf\{\tau \geq t \mid V_\tau \leq \tilde{D}_\tau\} \quad (6.72)$$

If  $\gamma_t > T$ , no default occurs over the time interval in question.

Therefore, the investor wants to optimize

$$u(x, y, t) = \max_{\pi_s(\cdot)} E [U(Y_T)\mathbb{I}_{\{\gamma_t > T\}} + U(Y_{\gamma_t}e^{r(T-\gamma_t)})\mathbb{I}_{\{\gamma_t \leq T\}} \mid V_t = x, Y_t = y] \quad (6.73)$$

Here  $U(y)$  is a utility function ( $U(y) = -e^{-hy}$  is used in the paper,  $h > 0$ ). In this case, the HJB equation is

$$u_t + \mathcal{L}u + rxu_x + \max_{\pi} \left[ \frac{1}{2}\sigma^2\pi^2 u_{xx} + \pi(\rho\sigma\eta u_{xy} + (\mu - r)u_x) \right] = 0 \quad (6.74)$$

where  $\mathcal{L}$  is the operator

$$\mathcal{L}u = \frac{\eta^2 y^2}{2} u_{yy} + \nu y u_y. \quad (6.75)$$

The domain for  $u$  is  $\{(t, y, x) \mid t \in [0, T], x \in \mathbb{R}, y \in [\tilde{D}(t), \infty)\}$ , and the boundary condition is

$$u(x, y, T) = U(x), \quad u(x, D e^{-\beta(T-t)}, t) = U(x e^{r(T-t)}). \quad (6.76)$$

There are other examples in this paper illustrating techniques for valuation of credit derivatives.

## 6.8 The basic theory of Hamilton-Jacobi equations

### The Euler-Lagrange equations

We now describe the approach to the Hamilton-Jacobi equations in terms of calculus of variations rather than optimal control.

Let  $L(q, x)$  be a smooth function,  $q, x \in \mathbb{R}^n$  called the Lagrangian. Fix two points  $x, y \in \mathbb{R}^n$  and consider the class of admissible functions

$$\mathcal{A} = \{w \in C([0, t]; \mathbb{R}^n) : w(0) = y, w(t) = x\},$$

that is  $w(t)$  are smooth paths that start at  $y$  at time zero, and end at  $x$  at time  $t$ . Define the functional

$$I(w) = \int_0^t L(\dot{w}(s), w(s)) ds.$$

The basic problem of the calculus of variations is to find the optimal curve  $w(t)$ :

$$\text{find } I^* = \min_{w \in \mathcal{A}} I(w),$$

and, if possible, the optimal path  $z(s) \in \mathcal{A}$  such that  $I(z) = I^*$ . Let us first assume that such  $z(s)$  exists and deduce some of its properties.

**Theorem 6.4** (*Euler-Lagrange equations*) *The function  $z(s)$  satisfies the Euler-Lagrange equations*

$$-\frac{d}{ds}[\nabla_q L(\dot{z}(s), z(s))] + \nabla_x L(\dot{z}(s), z(s)) = 0, \quad 0 \leq s \leq t. \quad (6.77)$$

**Proof.** Let  $v(t)$  be a smooth function such that  $v(0) = v(t) = 0$  and consider  $w_\tau(s) = z(s) + \tau v(s)$ . Set also  $r(\tau) = I(w_\tau)$ . As  $z(s)$  minimizes  $I(w)$  over  $\mathcal{A}$  and  $w_\tau \in \mathcal{A}$  for all  $\tau$ , we have  $r'(0) = 0$ . Let us now compute  $r'(\tau)$ :

$$r(\tau) = \int_0^t L(\dot{z}(s) + \tau \dot{v}(s), z(s) + \tau v(s)) ds,$$

so

$$r'(\tau) = \int_0^t [\nabla_q L \cdot \dot{v}(s) + \nabla_x L \cdot v(s)] ds = \int_0^t \left[-\frac{d}{ds} \nabla_q L + \nabla_x L\right] \cdot v(s) ds.$$

We integrated by parts in the second equality above, and used the fact that the boundary terms vanish since  $v(0) = v(t) = 0$ . Since  $r'(0) = 0$  for all  $v(s)$  as above, we should have

$$-\frac{d}{ds} \nabla_q L(\dot{z}(s), z(s)) + \nabla_x L(\dot{z}(s), z(s)) = 0,$$

which is (6.77).  $\square$

The above computation shows that if  $z(s)$  is a minimizer then it has to satisfy the Euler-Lagrange equation (6.77). However, of course, it is possible that  $z(s)$  is a critical point of  $I(w)$  but not its minimum – in that case  $z(s)$  also satisfies the Euler-Lagrange equations.

## The Hamilton equations

There is a nice connection between the Euler-Lagrange equations and the Hamilton equations of classical mechanics. We assume that the equation

$$p = \nabla_q L(q, x) \quad (6.78)$$

can be solved uniquely as an equation for  $q$ , as a smooth function of  $p$  and  $x$ . If that is the case, we can define the Hamiltonian

$$H(p, x) = p \cdot q(p, x) - L(q(p, x), x), \quad (6.79)$$

with the function  $q(p, x)$  defined implicitly by (6.78).

Let us now assume that  $z(s)$  is the solution of the Euler-Lagrange equations, and set

$$p(s) = \nabla_q L(\dot{z}(s), z(s)), \quad (6.80)$$

that is,

$$\dot{z}(s) = q(p(s), z(s)). \quad (6.81)$$

Differentiating (6.79) in  $p_j$  gives

$$\frac{\partial H(p(s), z(s))}{\partial p_j} = q_j(p(s), z(s)) + \sum_{i=1}^n p_i(s) \frac{\partial q_i}{\partial p_j} - \sum_{i=1}^m \frac{\partial L}{\partial q_i} \frac{\partial q_i}{\partial p_j} = q_j.$$

We used (6.80) in the last step. Using this in (6.81) gives

$$\dot{z}_j(s) = \frac{\partial H(p(s), z(s))}{\partial p_j}. \quad (6.82)$$

The Euler-Lagrange equations say that

$$\dot{p}_j(s) = \frac{\partial L}{\partial x_j}. \quad (6.83)$$

Differentiating (6.79) in  $x$  gives:

$$\frac{\partial H}{\partial x_j} = \sum_{i=1}^n p_i \frac{\partial q_i}{\partial x_j} - \frac{\partial L}{\partial x_j} - \sum_{i=1}^n \frac{\partial L}{\partial q_i} \frac{\partial q_i}{\partial x_j} = -\frac{\partial L}{\partial x_j}.$$

Now, putting this together with (6.82)-(6.83) gives the Hamiltonian system

$$\dot{z}(s) = \nabla_p H(p(s), z(s)), \quad \dot{p}(s) = -\nabla_z H(p(s), z(s)). \quad (6.84)$$

### The Legendre transform

Let us now assume that the Lagrangian does not depend on the variable  $x$ :  $L = L(q)$ . Then the Hamiltonian  $H(p)$  is

$$H(p) = p \cdot q(p) - L(q(p)), \quad (6.85)$$

with  $q$  being the solution of

$$p = \nabla_q L(q). \quad (6.86)$$

In order to ensure that the function  $q(p)$  is well-defined let us assume that the function  $L(q)$  is convex and

$$\lim_{|q| \rightarrow +\infty} \frac{L(q)}{|q|} = +\infty. \quad (6.87)$$

Let us now fix  $p$  and consider the function  $r(q) = p \cdot q - L(q)$ . This function is convex and  $r(q) \rightarrow -\infty$  as  $|q| \rightarrow +\infty$ . Therefore,  $r(q)$  attains a unique maximum at the point where  $p = \nabla L(q)$ , which is exactly the same equation as (6.86). Therefore, we may reformulate (6.85) as

$$H(p) = \sup_q (p \cdot q - L(q)). \quad (6.88)$$

The function  $H(p)$  defined by (6.88) is called the Legendre transform of  $L(q)$ , denoted as  $H(p) = L^*(q)$ .

**Theorem 6.5** *Assume that the function  $L(q)$  is convex and (6.87) holds, then  $H(p)$  is also convex, and*

$$\lim_{|p| \rightarrow +\infty} \frac{H(p)}{|p|} = +\infty. \quad (6.89)$$

*Moreover,  $L(q)$  is the Legendre transform of the function  $H$ .*

**Proof.** The function  $s(p; q) = p \cdot q - L(q)$  is an affine function of  $p$  for each  $q$  fixed. Therefore,  $H(p)$  is a supremum of a family of affine functions – hence, it is convex. Indeed, for any  $\lambda \in (0, 1)$  we have

$$\begin{aligned} H(\lambda p_1 + (1 - \lambda)p_2) &= \sup_q (\lambda p_1 + (1 - \lambda)p_2 \cdot q) - L(q) \\ &= \sup_q [\lambda p_1 \cdot q - \lambda L(q) + (1 - \lambda)p_2 \cdot q - (1 - \lambda)L(q)] \\ &\leq \sup_q [\lambda p_1 \cdot q - \lambda L(q)] + \sup_q [(1 - \lambda)p_2 \cdot q - (1 - \lambda)L(q)] = \lambda H(p_1) + (1 - \lambda)H(p_2), \end{aligned}$$

hence  $H(p)$  is convex.

In order to see that (6.89) holds, fix  $\lambda > 0$  and take  $\bar{q} = \lambda p / |p|$  in the definition of  $H(p)$ , then  $|\bar{q}| \leq \lambda$ , hence

$$H(p) \geq p \cdot \bar{q} - L(\bar{q}) = \lambda |p| - L(\bar{q}) \geq \lambda |p| - \sup_{|q| \leq \lambda} L(q).$$

It follows that

$$\lim_{|p| \rightarrow +\infty} \frac{H(p)}{|p|} \geq \lambda,$$

for each  $\lambda > 0$ , thus (6.89) holds.

In order to show that  $L(q)$  is actually the Legendre transform of  $H(p)$ , note that, for all  $p$  and  $q$  we have

$$H(p) \geq p \cdot q - L(q),$$

whence

$$L(q) \geq p \cdot q - H(p).$$

It follows that  $L(q) \geq H^*(q)$ . But we also have

$$H^*(q) = \sup_p [p \cdot q - H(p)] = \sup_p [p \cdot q - \sup_y [p \cdot y - L(y)]] = \sup_p \inf_y [p \cdot (q - y) + L(y)]. \quad (6.90)$$

As the function  $L(q)$  is convex, for each  $q$  there exists  $s(q)$  such that the graph of  $L(y)$  lies above the corresponding hyperplane:

$$L(y) \geq L(q) + s \cdot (y - q).$$

Let us take  $p = s(q)$  in (6.90):

$$H^*(q) \geq \inf_y [s \cdot (q - y) + L(y)] \geq L(q). \quad (6.91)$$

We conclude that  $H^*(p) = L(q)$ .  $\square$

## The Hopf-Lax formula

We now relate the variational problem that looked at to the Hamilton-Jacobi equations. Consider the initial value problem

$$u_t + H(\nabla u) = 0, \quad t > 0, \quad x \in \mathbb{R}^n, \quad (6.92)$$

with the initial data  $u(0, x) = g(x)$ . The initial data  $g(x)$  is globally Lipschitz continuous:

$$\text{Lip}(g) = \sup_{x, y \in \mathbb{R}^n} \frac{|g(x) - g(y)|}{|x - y|} < +\infty. \quad (6.93)$$

We assume that  $H(p)$  is convex and satisfies the growth condition (6.89). Let us define

$$u(t, x) = \inf \left[ \int_0^t L(\dot{w}(s)) ds + g(y) : w(0) = y, w(t) = x \right], \quad (6.94)$$

with the infimum taken over all  $C^1$  functions  $w(t)$  that satisfy the constraint  $w(t) = x$ . Here  $L(q)$  is the Legendre transform of the function  $H(p)$ . We will show that expression (6.94) gives a solution of the Hamilton-Jacobi equation (6.92).

**Theorem 6.6** (*Hopf-Lax formula*) *The function  $u(t, x)$  defined by (6.94) can be written as*

$$u(t, x) = \min_{y \in \mathbb{R}^n} \left[ tL \left( \frac{x - y}{t} \right) + g(y) \right]. \quad (6.95)$$

**Proof.** First, for any  $y \in \mathbb{R}^n$  we may take a "test path"

$$w(s) = y + \frac{s}{t}(x - y),$$

leading to

$$u(t, x) \leq \int_0^t L\left(\frac{x - y}{t}\right) ds + g(y) = tL\left(\frac{x - y}{t}\right) + g(y).$$

As a consequence, we have

$$u(t, x) \leq \inf_{y \in \mathbb{R}^n} \left[ tL\left(\frac{x - y}{t}\right) + g(y) \right].$$

On the other hand, Jensen's inequality implies that for any test path  $w(s)$  we have

$$\frac{1}{t} \int_0^t L(\dot{w}(s)) ds \geq L\left(\frac{1}{t} \int_0^t \dot{w}(s) ds\right).$$

Therefore,

$$\int_0^t L(\dot{w}(s)) ds \geq tL\left(\frac{x - y}{t}\right),$$

where  $y = w(0)$ , and thus

$$u(t, x) \geq \inf_{y \in \mathbb{R}^n} \left[ tL\left(\frac{x - y}{t}\right) + g(y) \right].$$

Thus, we have shown that

$$u(t, x) = \inf_{y \in \mathbb{R}^n} \left[ tL\left(\frac{x - y}{t}\right) + g(y) \right].$$

The fact that the infimum in the right side is actually achieved follow from the fact that for each  $t$  and  $x$  fixed the function

$$r(y) = tL\left(\frac{x - y}{t}\right) + g(y)$$

tends to  $+\infty$  as  $|y| \rightarrow +\infty$ . This is because  $L(y)$  is super-linear at infinity, and  $g$  is globally Lipschitz.  $\square$

## A formal computation of the Hamilton-Jacobi equation

Let us now show why we expect the function given by the Hopf-Lax formula to satisfy the hamilton-Jacobi equation, assuming that it is as smooth as needed. For simplicity, assume that  $x \in \mathbb{R}$ . Let  $z$  be such that

$$u(t, x) = tL\left(\frac{x-z}{t}\right) + g(z).$$

Then  $z$  is determined by the condition

$$g'(z) = L'\left(\frac{x-z}{t}\right), \quad (6.96)$$

hence we have

$$u_t = L\left(\frac{x-z}{t}\right) - L'\left(\frac{x-z}{t}\right)z_t - \frac{(x-z)}{t}L'\left(\frac{x-z}{t}\right) + g'(z)z_t = L\left(\frac{x-z}{t}\right) - \frac{(x-z)}{t}L'\left(\frac{x-z}{t}\right).$$

Moreover,

$$u_x = L'\left(\frac{x-z}{t}\right), \quad (6.97)$$

hence the above can be written as

$$u_t = L\left(\frac{x-z}{t}\right) - u_x \frac{(x-z)}{t}. \quad (6.98)$$

On the other hand, in the definition of  $H(p)$  we have

$$H(p) = \sup_{y \in \mathbb{R}} (py - L(y)) = pq - L(q),$$

with  $q$  determined by the relation  $p = L'(q)$ . Therefore,

$$H(u_x) = u_x q - L(q),$$

with  $q$  such that  $u_x = L'(q)$ . But (6.97) implies that then  $q = (x-z)/t$ , and (6.98) is nothing but the Hamilton-Jacobi equation

$$u_t + H(u_x) = 0.$$

## The rigorous derivation of the Hamilton-Jacobi equation

Let us now verify that the Hopf-Lax formula is Lipschitz continuous.

**Lemma 6.7** *Let  $u(t, x)$  be defined by (6.95). Then the function  $u(t, x)$  is Lipschitz continuous in  $x$  for  $t \geq 0$  and  $x \in \mathbb{R}^n$ , and  $u(t, x) \rightarrow g(x)$  as  $t \rightarrow 0$ .*

**Proof.** Take  $x_1, x_2 \in \mathbb{R}^n$ , and choose  $y$  so that

$$u(t, x_1) = tL\left(\frac{x_1 - y}{t}\right) + g(y),$$

then, choosing  $z = x_2 - x_1 + y$  below, gives

$$\begin{aligned} u(t, x_1) - u(t, x_2) &= \min_{z \in \mathbb{R}^n} \left[ tL \left( \frac{x_2 - z}{t} \right) + g(z) \right] - tL \left( \frac{x_1 - y}{t} \right) - g(y) \\ &\leq g(x_2 - x_1 + y) - g(y) \leq \text{Lip}(g) |x_1 - x_2|. \end{aligned}$$

Switching the roles of  $x_1$  and  $x_2$  gives Lipschitz continuity in  $x$ :

$$|u(t, x_1) - u(t, x_2)| \leq \text{Lip}(g) |x_1 - x_2|.$$

In order to verify the initial condition, note that choosing  $y = x$  gives

$$u(t, x) \leq tL(0) + g(x), \tag{6.99}$$

but we also have

$$\begin{aligned} u(t, x) &= \min_y \left[ tL \left( \frac{x - y}{t} \right) + g(y) \right] \geq \min_y \left[ tL \left( \frac{x - y}{t} \right) + g(x) - \text{Lip}(g) |x - y| \right] \\ &= g(x) + \min_z [tL(z) - \text{Lip}(g)t|z|] = g(x) + t \min_z [L(z) - \text{Lip}(g)|z|]. \end{aligned}$$

once again, as  $L(z)$  grows super-linearly at infinity, we have

$$\min_z [L(z) - \text{Lip}(g)|z|] > -\infty,$$

hence

$$u(t, x) \geq g(x) - Ct. \tag{6.100}$$

We conclude that  $u(t, x) \rightarrow g(x)$  as  $t \rightarrow 0$ .  $\square$

In order to show that  $u(t, x)$  is Lipschitz continuous in time, we need the following lemma (which is essentially a version of the dynamic programming principle).

**Lemma 6.8** *For each  $x \in \mathbb{R}^n$ , and  $0 \leq s < t$  we have*

$$u(t, x) = \min_{y \in \mathbb{R}^n} \left[ (t - s)L \left( \frac{x - y}{t - s} \right) + u(s, y) \right]. \tag{6.101}$$

**Proof.** Choose  $z$  so that

$$u(s, y) = sL \left( \frac{y - z}{s} \right) + g(z).$$

Let us write

$$\frac{x - z}{t} = \left(1 - \frac{s}{t}\right) \frac{x - y}{t - s} + \frac{s}{t} \frac{y - z}{s}.$$

As  $L$  is convex, it follows that

$$\begin{aligned} u(t, x) &\leq tL \left( \frac{x - z}{t} \right) + g(z) \leq t \left(1 - \frac{s}{t}\right) L \left( \frac{x - y}{t - s} \right) + sL \left( \frac{y - z}{s} \right) + g(z) \\ &= (t - s)L \left( \frac{x - y}{t - s} \right) + sL \left( \frac{y - z}{s} \right) + g(z) = (t - s)L \left( \frac{x - y}{t - s} \right) + u(s, y), \end{aligned}$$

and thus

$$u(t, x) \leq \inf_{y \in \mathbb{R}^n} \left[ (t-s)L\left(\frac{x-y}{t-s}\right) + u(s, y) \right].$$

As the function  $u(s, y)$  is actually continuous in  $y$  (this follows from Lemma 6.7), and  $|u(s, y)|$  grows not faster than linearly at infinity (that follows from (6.99)-(6.100)), the infimum in the right side is actually attained:

$$u(t, x) \leq \min_{y \in \mathbb{R}^n} \left[ (t-s)L\left(\frac{x-y}{t-s}\right) + u(s, y) \right].$$

In order to show the opposite inequality, choose  $z$  so that

$$u(t, x) = tL\left(\frac{x-z}{t}\right) + g(z),$$

and set

$$y = \frac{s}{t}x + \left(1 - \frac{s}{t}\right)z.$$

Then, we have

$$\frac{x-y}{t-s} = \frac{x-z}{t} = \frac{y-z}{s},$$

hence

$$(t-s)L\left(\frac{x-y}{t-s}\right) + u(s, y) \leq (t-s)L\left(\frac{x-z}{t}\right) + sL\left(\frac{y-z}{s}\right) + g(z) = tL\left(\frac{x-z}{t}\right) + g(z) = u(t, x).$$

This proves (6.101).  $\square$

**Lemma 6.9** *The function  $u(t, x)$  defined by (6.95) is Lipschitz continuous in  $t$  for  $t \geq 0$  and  $x \in \mathbb{R}^n$ .*

**Proof.** Combing the ideas in the proof of Lemma 6.7 (see (6.99)-(6.100)) with the result of Lemma 6.8 gives

$$u(s, x) - C(t-s) \leq u(t, x) \leq u(s, x) + C(t-s),$$

and we are done.  $\square$

Since the function  $u(t, x)$  is Lipschitz in  $t$  and  $x$ , it is differentiable almost everywhere.

**Theorem 6.10** *The function  $u(t, x)$  defined by (6.95) is Lipschitz continuous in  $t$  and  $x$ , differentiable almost everywhere and solves the initial value problem*

$$u_t + H(\nabla u) = 0, \quad t > 0, \quad x \in \mathbb{R}^n, \quad (6.102)$$

with  $u(0, x) = g(x)$ .

**Proof.** It remains only to verify that at the points  $(t, x)$  where both  $u_t$  and  $\nabla u$  exist, the Hamilton-Jacobi equation (6.102) is satisfied. Fix  $q \in \mathbb{R}^n$ ,  $h > 0$ , then we have, according to Lemma 6.8:

$$u(x + hq, t + h) = \min_{y \in \mathbb{R}^n} \left[ hL\left(\frac{x + hq - y}{h}\right) + u(t, y) \right] \leq hL(q) + u(t, x).$$

It follows that

$$u_t(t, x) + q \cdot \nabla u(t, x) \leq L(q),$$

for all  $q \in \mathbb{R}^n$ . Therefore, we have

$$u_t(t, x) + H(\nabla u(t, x)) = u_t(t, x) + \min_{q \in \mathbb{R}^n} (q \cdot \nabla u(t, x) - L(q)) \leq 0. \quad (6.103)$$

Next, we show the opposite inequality. Choose  $z$  so that

$$u(t, x) = tL\left(\frac{x-z}{t}\right) + g(z).$$

Given  $h > 0$ , set

$$y = \frac{t-h}{t}x + \left(1 - \frac{t-h}{t}\right)z = x - h\frac{(x-z)}{t}, \quad (6.104)$$

so that

$$\frac{x-z}{t} = \frac{y-z}{t-h}.$$

We have

$$u(t, x) - u(t-h, y) \geq tL\left(\frac{x-z}{t}\right) + g(z) - \left[(t-h)L\left(\frac{y-z}{t-h}\right) + g(z)\right] = hL\left(\frac{x-z}{t}\right).$$

Keeping in mind expression (6.104), and letting  $h \rightarrow 0$  gives

$$u_t(t, x) + \frac{1}{t}(x-z) \cdot \nabla u(t, x) \geq L\left(\frac{x-z}{t}\right).$$

It follows that

$$u_t(t, x) + H(\nabla u(t, x)) = u_t(t, x) + \min_{q \in \mathbb{R}^n} (q \cdot \nabla u(t, x) - L(q)) \geq 0,$$

which, together with (6.103) finishes the proof.  $\square$

## 6.9 Viscosity solutions for Hamilton-Jacobi equations

We will now consider solutions of the Cauchy problem for the Hamilton-Jacobi equations

$$\begin{aligned} u_t + H(\nabla u, x) &= 0, \quad t \geq 0, \quad x \in \mathbb{R}^n, \\ u(0, x) &= g(x). \end{aligned} \quad (6.105)$$

The idea is to consider solutions of the regularized parabolic problem

$$\begin{aligned} u_t^\varepsilon + H(\nabla u^\varepsilon, x) &= \varepsilon \Delta u^\varepsilon, \\ u^\varepsilon(0, x) &= g(x). \end{aligned} \quad (6.106)$$

The idea is to show that for each  $\varepsilon > 0$  the problem (6.105) admits a regular solution, and then pass to the limit  $\varepsilon \rightarrow 0$ . The difficulty is that as  $\varepsilon \rightarrow 0$  the regularizing effect of the Laplacian is less and less, so  $u^\varepsilon(t, x)$  are less and less regular, and it is not clear that the limit

of  $u^\varepsilon(t, x)$ , if it exists, is regular, and in which sense it would satisfy the Hamilton-Jacobi equation (6.106).

Let us for the moment assume that for some sequence  $\varepsilon_n \rightarrow 0$  (6.105) has a smooth solution  $u_n(t, x) = u^{\varepsilon_n}(t, x)$ , and that  $u_n(t, x) \rightarrow u(t, x)$  locally uniformly in  $\mathbb{R}^n \times [0, +\infty)$ . Let us take a smooth test function  $v$  and suppose that  $u(t, x) - v(t, x)$  has a strict local minimum at some point  $(t_0, x_0)$ . Then  $u(t, x) - v(t, x) > u(t_0, x_0) - v(t_0, x_0)$  in some neighborhood  $B$  of  $(t_0, x_0)$ . The functions  $u_n(t, x) - v(t, x)$  have to attain a local maximum inside  $B$  when  $n$  is sufficiently large as well – simply because we have

$$\max_{\partial B} (u_n(t, x) - v(t, x)) < u_n(t_0, x_0) - v(t_0, x_0).$$

Hence,  $u_n(t, x) - v(t, x)$  attains a maximum in  $B$ . Now, if we let the radius of  $B$  go to zero, we get a sequence of points  $(t_n, x_n) \rightarrow (t_0, x_0)$  such that  $u_n(t, x) - v(t, x)$  has a local maximum at  $(t_n, x_n)$ . We deduce that  $\nabla u_n(t_n, x_n) = \nabla v(t_n, x_n)$ ,  $u_{n,t}(t_n, x_n) = v_t(t_n, x_n)$  and

$$-\Delta u_n(t_n, x_n) \geq -\Delta v(t_n, x_n).$$

It follows that

$$\begin{aligned} v_t(t_n, x_n) + H(\nabla v(t_n, x_n), x_n) &= u_{n,t}(t_n, x_n) + H(\nabla u_n(t_n, x_n), x_n) = \varepsilon \Delta u_n(t_n, x_n) \\ &\leq \varepsilon \Delta v(t_n, x_n). \end{aligned} \tag{6.107}$$

The function  $v(t, x)$  is smooth, so we may let  $\varepsilon \rightarrow 0$  in (6.107) to conclude that

$$v_t(t_0, x_0) + H(\nabla v(t_0, x_0), x_0) \leq 0. \tag{6.108}$$

Inequality (6.108) should hold for any smooth function  $v(t, x)$  such that  $u(t, x) - v(t, x)$  attains a local maximum at  $(t_0, x_0)$ . Similarly, if  $u - v$  attains a local minimum at  $(t_0, x_0)$  then we should have

$$v_t(t_0, x_0) + H(\nabla v(t_0, x_0), x_0) \geq 0. \tag{6.109}$$

The above argument assumed that solutions of the regularized parabolic problem exist and have a limit  $u(t, x)$ . Let us now instead take the inequalities (6.108) and (6.109) as the starting point and define the appropriate solution of the Hamilton-Jacobi equation purely in their terms, forgetting everything about the parabolic problem.

**Definition 6.11** *A bounded uniformly continuous function  $u(t, x)$  is a viscosity solution of the Cauchy problem (6.105) for the Hamilton-Jacobi equation if  $u(0, x) = g(x)$  for all  $x \in \mathbb{R}^n$ , and for each  $v \in C^\infty([0, \infty) \times \mathbb{R}^n)$  such that  $u - v$  has a local maximum at a point  $(t_0, x_0)$  with  $t_0 > 0$ , we have*

$$v_t(t_0, x_0) + H(\nabla v(t_0, x_0), x_n) \leq 0, \tag{6.110}$$

*while if  $u - v$  attains a local minimum at a point  $(t_0, x_0)$  with  $t_0 > 0$ , we have*

$$v_t(t_0, x_0) + H(\nabla v(t_0, x_0), x_n) \geq 0, \tag{6.111}$$

We will verify that these two conditions are reasonable in the following sense.

**Theorem 6.12** (Consistency) *Let  $u(t, x)$  be a viscosity solution of the Cauchy problem*

$$\begin{aligned} u_t + H(\nabla u, x) &= 0, \quad t \geq 0, \quad x \in \mathbb{R}^n, \\ u(0, x) &= g(x). \end{aligned} \tag{6.112}$$

*Assume that  $u$  is differentiable at some point  $(t_0, x_0)$  with  $t_0 > 0$ , then*

$$u_t(t_0, x_0) + H(\nabla u(t_0, x_0), x_0) = 0. \tag{6.113}$$

We begin the proof with the following lemma.

**Lemma 6.13** *Assume that  $u(x)$ ,  $x \in \mathbb{R}^n$  is a continuous function and  $u(x)$  is differentiable at  $x_0$ . Then there exists a  $C^1(\mathbb{R}^n)$  function  $q(x)$  such that  $u(x_0) = v(x_0)$  and  $u - v$  has a strict local maximum at  $x_0$ .*

**Proof of Lemma.** Let us set

$$v(x) = u(x + x_0) - u(x_0) - x \cdot \nabla u(x_0),$$

so that  $v(0) = 0$ ,  $\nabla v(0) = 0$ . It follows that  $v(x) = |x|\rho(x)$ , where the function  $\rho(x)$  is continuous, and  $\rho(0) = 0$ . Set

$$p(r) = \max_{x \in B(0, r)} \rho(x),$$

then  $p(r)$  is continuous, non-decreasing and  $p(0) = 0$ . Finally, define

$$w(x) = |x|^2 + \int_{|x|}^{2|x|} p(r) dr.$$

Then  $w \in C^1(\mathbb{R}^n)$ , and

$$|w(x)| \leq |x|^2 + |x|p(2|x|),$$

which means that  $w(0) = 0$  and  $\nabla w(0) = 0$ . However, we have

$$\begin{aligned} v(x) - w(x) &= |x|\rho(x) - |x|^2 - \int_{|x|}^{2|x|} p(r) dr \leq |x|p(|x|) - |x|^2 - \int_{|x|}^{2|x|} p(r) dr \\ &\leq -|x|^2 < 0 = v(0) - w(0). \end{aligned}$$

Therefore, the function  $v(x) - w(x)$  attains its local maximum at  $x = 0$ , which means that we can take

$$q(x) = w(x - x_0) + u(x_0) + (x - x_0) \cdot \nabla u(x_0),$$

proving Lemma 6.13.  $\square$

**Proof of Theorem 6.14.** Note that if  $u(t, x)$  were  $C^\infty$  (rather than just differentiable at  $(t_0, x_0)$ ), we could take  $u$  itself as a test function in the definition of the viscosity solution, and conclude that both hence

$$u_t(t_0, x_0) + H(\nabla u(t_0, x_0), x_0) \leq 0,$$

and

$$u_t(t_0, x_0) + H(\nabla u(t_0, x_0), x_0) \geq 0,$$

giving the result. Hence, what we need to do is replace  $u$  by a smooth test function without changing  $u_t$  and  $\nabla u$  too much. Lemma 6.13 implies that there exists a  $C^1$  function  $v$  such that  $u - v$  has a strict maximum at  $(x_0, t_0)$ . Next, let  $v_\varepsilon(t, x)$  be

$$v_\varepsilon(t, x) = \frac{1}{\varepsilon^{n+1}} \int \chi\left(\frac{t-s}{\varepsilon}, \frac{x-y}{\varepsilon}\right) v(s, y) ds dy.$$

Here the function  $\chi(t, x) \in C^\infty(\mathbb{R}^{n+1})$  is chosen so that  $\chi(t, x) \geq 0$ , and

$$\int \chi(t, x) dt dx = 1.$$

Then the functions  $v_\varepsilon \in C^\infty$  for all  $\varepsilon > 0$ , and  $v_\varepsilon \rightarrow v$ ,  $v_{\varepsilon,t} \rightarrow v_t$ ,  $\nabla v_\varepsilon \rightarrow \nabla v$ , all locally uniformly near  $(t_0, x_0)$ . It follows that  $u(t, x) - v_\varepsilon(t, x)$  has a strict local maximum at some point  $(t_\varepsilon, x_\varepsilon)$  with  $(t_\varepsilon, x_\varepsilon) \rightarrow (t_0, x_0)$  as  $\varepsilon \rightarrow 0$ . The definition of the viscosity solution implies that

$$v_{\varepsilon,t}(t_\varepsilon, x_\varepsilon) + H(\nabla v_\varepsilon(t_\varepsilon, x_\varepsilon), x_\varepsilon) \leq 0.$$

Passing to the limit  $\varepsilon \rightarrow 0$  gives

$$v_t(t_0, x_0) + H(\nabla v(t_0, x_0), x_0) \leq 0.$$

Since  $u(t, x)$  is differentiable at  $(t_0, x_0)$  and  $u - v$  attains a local maximum at  $(t_0, x_0)$ , we have

$$u_t(t_0, x_0) = v_t(t_0, x_0), \quad \nabla u(t_0, x_0) = \nabla v(t_0, x_0),$$

hence

$$u_t(t_0, x_0) + H(\nabla u(t_0, x_0), x_0) \leq 0.$$

Similarly, we can prove that

$$v_t(t_0, x_0) + H(\nabla v(t_0, x_0), x_0) = 0,$$

and we are done.  $\square$

Viscosity solution (if it exists) is unique.

**Theorem 6.14** (*Uniqueness*) *There exists at most one viscosity solution of the Cauchy problem*

$$\begin{aligned} u_t + H(\nabla u, x) &= 0, \quad t \geq 0, \quad x \in \mathbb{R}^n, \\ u(0, x) &= g(x). \end{aligned} \tag{6.114}$$

### Hopf-Lax formula as a viscosity solution

Let us now show that the Hopf-Lax formula gives a viscosity solution for the Cauchy problem

$$\begin{aligned} u_t + H(\nabla u) &= 0, \quad t \geq 0, \quad x \in \mathbb{R}^n, \\ u(0, x) &= g(x), \end{aligned} \tag{6.115}$$

if  $H(p)$  is convex,

$$\lim_{|p| \rightarrow +\infty} \frac{H(p)}{|p|} = +\infty,$$

and  $g(x)$  is bounded and Lipschitz continuous. Let  $L$  be the Legendre transform of  $H$ :

$$L(q) = \sup_{p \in \mathbb{R}^n} (p \cdot q - H(p)),$$

and set

$$u(t, x) = \min_{y \in \mathbb{R}^n} \left[ tL\left(\frac{x-y}{t}\right) + g(y) \right]. \quad (6.116)$$

Let us show that  $u(t, x)$  is the viscosity solution of (6.115). We already know that  $u(t, x)$  defined by (6.116) is Lipschitz continuous in  $t$  and  $x$ .

Take  $v \in C^\infty$  and assume that  $u - v$  has a local maximum at  $(t_0, x_0)$ . Then, we have

$$u(t_0, x_0) = \min_{x \in \mathbb{R}^n} \left[ (t_0 - t)L\left(\frac{x_0 - x}{t_0 - t}\right) + u(t, x) \right] \leq (t_0 - t)L\left(\frac{x_0 - x}{t_0 - t}\right) + u(t, x),$$

for all  $0 \leq t < t_0$ , and  $x \in \mathbb{R}^n$ . Since  $u - v$  has a local maximum at  $(t_0, x_0)$ , we also have

$$u(t, x) - v(t, x) \leq u(t_0, x_0) - v(t_0, x_0),$$

for  $t, x$  close to  $t_0, x_0$ . Hence,

$$v(t_0, x_0) - v(t, x) \leq u(t_0, x_0) - u(t, x) \leq (t_0 - t)L\left(\frac{x_0 - x}{t_0 - t}\right).$$

Let us use this relation for  $t = t_0 - h$  and  $x = x_0 - hq$ , with some  $h > 0$  fixed, and  $q \in \mathbb{R}^n$ . We get

$$v(t_0, x_0) - v(t_0 - h, x_0 - hq) \leq hL(q).$$

Passing to the limit  $h \rightarrow 0$  gives

$$v_t + q \cdot \nabla v(t_0, x_0) \leq L(q).$$

As this is true for all  $q$ , we deduce that

$$v_t(t_0, x_0) + H(\nabla v(t_0, x_0)) \leq 0.$$

Next, suppose that  $u - v$  attains a local minimum at  $(t_0, x_0)$ . We will show that

$$v_t(t_0, x_0) + H(\nabla v(t_0, x_0)) \geq 0.$$

If this is false, then there exists some  $\theta > 0$  so that

$$v_t(t, x) + H(\nabla v(t, x)) \leq -\theta < 0,$$

for all  $t$  and  $x$  close to  $t_0, x_0$ . It follows that

$$v_t(t, x) + q \cdot \nabla v(t, x) - L(q) \leq -\theta, \quad (6.117)$$

for all such  $t, x$  and all  $q \in \mathbb{R}^n$ . Now, for  $h > 0$  small enough there exists  $x_1$  close to  $x_0$  so that

$$u(t_0, x_0) = hL\left(\frac{x_0 - x_1}{h}\right) + u(t_0 - h, x_1).$$

Let us look at (6.117) with  $q = (x_0 - x_1)/h$ , then we get

$$v(x_0, t_0) - v(t - h, x_1) \leq h \left( L\left(\frac{x_0 - x_1}{h}\right) - \theta \right).$$

But that means

$$v(x_0, t_0) - v(t - h, x_1) < u(x_0, t_0) - u(t - h, x_1).$$

This is a contradiction to  $u - v$  attaining a local minimum at  $(t_0, x_0)$ .  $\square$