



Renormalised Steepest Descent in Hilbert Space Converges to a Two-Point Attractor

LUC PRONZATO¹, HENRY P. WYNN² and ANATOLY A. ZHIGLJAVSKY³

¹*Laboratoire I3S, CNRS/Université de Nice-Sophia Antipolis, bât. Euclide, Les Algorithmes,
2000 route des Lucioles, BP 121, 06903 Sophia-Antipolis Cedex, France*

²*Department of Statistics, University of Warwick, Coventry CV4 7AL, U.K.*

³*School of Mathematics, Cardiff University, Senghennydd Road, Cardiff CF24 4YH, Wales, U.K.*

(Received: 19 August 1998; in final form: 7 December 2000)

Abstract. The result that for quadratic functions the classical steepest descent algorithm in \mathbb{R}^d converges locally to a two-point attractor was proved by Akaike. In this paper this result is proved for bounded quadratic operators in Hilbert space. The asymptotic rate of convergence is shown to depend on the starting point while, as expected, confirming the Kantorovich bounds. The introduction of a relaxation coefficient in the steepest-descent algorithm completely changes its behaviour, which may become chaotic. Different attractors are presented. We show that relaxation allows a significantly improved rate of convergence.

Mathematics Subject Classifications (2000): 90C25, 68Q25, 34D45.

Key words: asymptotic behaviour, steepest descent, gradient, Hilbert space, quadratic operator, period-2 cycles.

1. Introduction

In the recent book [7] the authors study the behaviour of optimization and search algorithms when at each iteration they are renormalized to a standard region. In Chapter 7 of the book the classical steepest descent is studied when the normalization is the unit sphere. Thus if x_k is the k th iterate, then one looks at $z_k = x_k / \|x_k\|$. The authors reproved a classical result of Akaike [1] which says that when the objective function is quadratic: $x^T A x - x^T y$ (with y being a fixed vector) z_k converges to a two-point attractor which lies in the space spanned by the eigenvectors corresponding to the smallest and largest eigenvalues of the matrix A .

In this paper the authors generalise this result to bounded quadratic operators in Hilbert space. The proof stems from the proof for \mathbb{R}^d but is considerably more technical. In both cases, as in [1], the method consists of converting the problem to a one containing a special type of operator on one-dimensional measures. The additional technicalities arise from the fact that in the Hilbert space case the measure may be continuous and arises from the spectral measure of the operator.

The steepest descent method is widely used in functional and infinite-dimensional numerical analysis as being convenient to implement, although, of course, in \mathbb{R}^d there are many preferable algorithms. Examples of its use are the solution of linear equations, solutions of partial differential equations and general ad hoc optimisation. An outstanding introduction is Chapter 15 in [4] where the bounds on convergence rates (named after Kantorovich) are discussed. An important aspect of the current result, as indeed for the \mathbb{R}^d result, is that the actual asymptotic rate of convergence, although satisfying Kantorovich bounds, depends on the starting point x_0 and is difficult to predict. This complex behaviour has consequences for stability which is discussed following the main results.

2. Asymptotic Behaviour of the Steepest Descent Algorithm

Let A be a bounded self-adjoint operator in a real Hilbert space \mathcal{H} with inner product (x, y) . Assume that A is positive, bounded below, and denote its boundaries by m and M ,

$$m = \inf_{\|x\|=1} (Ax, x), \quad M = \sup_{\|x\|=1} (Ax, x),$$

with $0 < m < M < \infty$. We apply the steepest descent method for minimizing the quadratic form

$$f(x) = \frac{1}{2}(Ax, x) - (x, y). \quad (1)$$

It is minimum at $x^* = A^{-1}y$, its directional derivative at x in the direction u is

$$\frac{\partial f(x)}{\partial u} = \frac{(Ax - y, u)}{\|u\|}.$$

The direction of steepest descent at x is $-g$, with $g = g(x)$ the gradient at x , namely $g = Ax - y$. The minimum of f in this direction is obtained for the optimum step-length

$$\alpha = \frac{(g, g)}{(Ag, g)}.$$

One iteration of the steepest descent algorithm is thus

$$x_{k+1} = x_k - \frac{(g_k, g_k)}{(Ag_k, g_k)} g_k, \quad (2)$$

with $g_k = Ax_k - y$ and x_0 some initial element in \mathcal{H} . Since $y = Ax^*$, we can rewrite the iteration above as

$$(x_{k+1} - x^*) = (x_k - x^*) - \frac{(g_k, g_k)}{(Ag_k, g_k)} g_k,$$

with $g_k = A(x_k - x^*)$, so that

$$g_{k+1} = g_k - \frac{(g_k, g_k)}{(Ag_k, g_k)} Ag_k.$$

Define the renormalised variable

$$z_k = \frac{g_k}{(g_k, g_k)^{1/2}}, \quad (3)$$

so that $(z_k, z_k) = 1$, $g_{k+1} = g_k - Ag_k/(Az_k, z_k)$ and

$$(g_{k+1}, g_{k+1}) = (g_k, g_k) \left[\frac{(A^2 z_k, z_k)}{(Az_k, z_k)^2} - 1 \right],$$

which gives

$$z_{k+1} = \frac{[(Az_k, z_k)I - A]z_k}{[(A^2 z_k, z_k) - (Az_k, z_k)^2]^{1/2}}. \quad (4)$$

In the special case where $\mathcal{H} = \mathbb{R}^d$, we can assume that A is already diagonalised, with eigenvalues $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_d$. The iteration (4) gives

$$[z_{k+1}]_i^2 = \frac{(\sum_{j=1}^d \lambda_j [z_k]_j^2 - \lambda_i)^2}{\sum_{j=1}^d \lambda_j^2 [z_k]_j^2 - (\sum_{j=1}^d \lambda_j [z_k]_j^2)^2} [z_k]_i^2, \quad (5)$$

with $[z_k]_i$ the i th component of z_k . We can then consider $[z_k]_i^2$ as a weight on the eigenvalue λ_i , so that (5) defines a transformation on discrete probability measures supported on $\{\lambda_1, \dots, \lambda_d\}$. The asymptotic behaviour of the sequence (z_k) generated by (5) is studied in [1, 2] and in Chapter 7 of [7]. The main result is that, assuming $0 < \lambda_1 < \lambda_2 \leq \dots \leq \lambda_{d-1} < \lambda_d$, the sequence (z_k) converges to a two-dimensional plane, spanned by the eigenvectors e_1, e_d associated with λ_1 and λ_d . More precisely, the attraction property can be stated as follows: choose z_0 such that $(z_0, e_1) > 0$, $(z_0, e_d) > 0$, then

$$z_{2k} \rightarrow \sqrt{p}e_1 + \sqrt{1-p}e_d, \quad z_{2k+1} \rightarrow \sqrt{1-p}e_1 - \sqrt{p}e_d$$

when $k \rightarrow \infty$,

where p is some number in $(0, 1)$, see Section 4 concerning the range of possible values for p .

This property has important consequences for the asymptotic rate of convergence of the steepest descent algorithm, see Section 5. However, the proofs used in the references above do not apply when \mathcal{H} is infinite-dimensional, and we present below a more general proof of this attraction theorem.

THEOREM 1. *Let A be a bounded self-adjoint operator in a Hilbert space \mathcal{H} , positive, with bounds m and M such that $0 < m < M < \infty$, and apply the*

steepest descent algorithm (2) for the minimisation of $f(x)$ given by (1), initialised at x_0 such that, for any ϵ , $0 < \epsilon < (M - m)/2$,

$$(E_{m+\epsilon} z_0, z_0) > 0 \quad \text{and} \quad (E_{M-\epsilon} z_0, z_0) < 1, \quad (6)$$

with $z_0 = z(x_0)$, see (3), and (E_λ) the spectral family of projections associated with A . The asymptotic behaviour of the renormalised gradient $z_k = z(x_k)$ is such that

$$z_{2k} = \sqrt{p} u_{2k} + \sqrt{1-p} v_{2k}, \quad z_{2k+1} = \sqrt{1-p} u_{2k+1} - \sqrt{p} v_{2k+1}, \quad (7)$$

with $\|u_n\| = \|v_n\| = 1 \forall n$, $\|Au_n - mu_n\| \rightarrow 0$, $\|Av_n - Mv_n\| \rightarrow 0$ as $n \rightarrow \infty$, and p some number in $(0, 1)$, depending on z_0 .

Proof. Since A is self-adjoint, its spectrum \mathcal{S}_A is a closed subset of the interval $[m, M]$ of the real line and $m, M \in \mathcal{S}_A$. Let E_λ be the spectral family associated with A , and define the measure $\nu_k = d(E_\lambda z_k, z_k)$, $m \leq \lambda \leq M$. Since $(z_k, z_k) = \int_m^M \nu_k(d\lambda) = 1$, ν_k is a probability measure on the Borel sets of $(0, \infty)$, with $\nu_k([m, M]) = 1 \forall k$. This representation gives

$$(Az_k, z_k) = \int \lambda \nu_k(d\lambda), \quad (A^2 z_k, z_k) = \int \lambda^2 \nu_k(d\lambda),$$

where integration is over $[m, M]$ unless otherwise specified. Therefore, the transformation (4) gives in terms of ν_k :

$$\nu_{k+1}(\mathcal{A}) = \frac{\int_{\mathcal{A}} [\lambda - \int \lambda' \nu_k(d\lambda')]^2 \nu_k(d\lambda)}{\int \lambda'^2 \nu_k(d\lambda') - [\int \lambda' \nu_k(d\lambda')]^2}.$$

The conditions (6) on z_0 are equivalent to $\text{ess inf}(\nu_0) = m$ and $\text{ess sup}(\nu_0) = M$, see Theorem 2, and the updating rule for ν_k can be written as (8). Theorem 2 then implies (9), which can be written as: $\forall \epsilon > 0$, $\epsilon \leq \beta = (M - m)/2$,

$$\begin{aligned} (E_{m+\epsilon} z_{2k}, z_{2k}) &\rightarrow p, & (E_{M-\epsilon} z_{2k}, z_{2k}) &\rightarrow p, \\ (E_{m+\epsilon} z_{2k+1}, z_{2k+1}) &\rightarrow 1-p, & (E_{M-\epsilon} z_{2k+1}, z_{2k+1}) &\rightarrow 1-p, \end{aligned}$$

as $k \rightarrow \infty$, where p depends on z_0 , $0 < p < 1$. Define $p_{2k} = (E_{m+\beta} z_{2k}, z_{2k})$, $p_{2k+1} = 1 - (E_{m+\beta} z_{2k+1}, z_{2k+1})$, and the angles φ , φ_n by $\cos \varphi = \sqrt{p}$, $\sin \varphi = \sqrt{1-p}$, $\cos \varphi_n = \sqrt{p_n}$, $\sin \varphi_n = \sqrt{1-p_n}$, $\forall n$. Also define $s_{2k} = E_{m+\beta} z_{2k} / \cos \varphi_{2k}$, $s_{2k+1} = E_{m+\beta} z_{2k+1} / \sin \varphi_{2k+1}$, $t_{2k} = (z_{2k} - E_{m+\beta} z_{2k}) / \sin \varphi_{2k}$, $t_{2k+1} = -(z_{2k+1} - E_{m+\beta} z_{2k+1}) / \cos \varphi_{2k+1}$. This gives $p_n \rightarrow p$ as $n \rightarrow \infty$, $\|s_n\| = \|t_n\| = 1 \forall n$, and $z_{2k} = \cos \varphi_{2k} s_{2k} + \sin \varphi_{2k} t_{2k}$, $z_{2k+1} = \sin \varphi_{2k+1} s_{2k+1} - \cos \varphi_{2k+1} t_{2k+1}$. Also,

$$\|As_n - ms_n\|^2 = \int (\lambda - m)^2 d(E\lambda s_n, s_n),$$

which, for $n = 2k$ and any $\epsilon, 0 < \epsilon < \beta$, gives

$$\begin{aligned} \|As_{2k} - ms_{2k}\|^2 &= \int_m^{m+\beta} \frac{(\lambda - m)^2}{p_{2k}} d(E_\lambda z_{2k}, z_{2k}) \\ &= \int_m^{m+\epsilon} \frac{(\lambda - m)^2}{p_{2k}} d(E_\lambda z_{2k}, z_{2k}) + \\ &\quad + \int_{m+\epsilon}^{m+\beta} \frac{(\lambda - m)^2}{p_{2k}} d(E_\lambda z_{2k}, z_{2k}) \\ &\leq \frac{\epsilon^2}{p_{2k}} + \frac{\beta^2}{p_{2k}} \left[p_{2k} - \int_m^{m+\epsilon} d(E_\lambda z_{2k}, z_{2k}) \right]. \end{aligned}$$

Since $p_{2k} \rightarrow p$ and $\int_m^{m+\epsilon} d(E_\lambda z_{2k}, z_{2k}) \rightarrow p$ as $k \rightarrow \infty$, $\|As_{2k} - ms_{2k}\| \rightarrow 0$ as $k \rightarrow \infty$. Similarly, $\|As_{2k+1} - ms_{2k+1}\| \rightarrow 0$ as $k \rightarrow \infty$ and $\|At_n - mt_n\| \rightarrow 0$ as $n \rightarrow \infty$. Consider now

$$u_n = \cos \vartheta_n s_n + \sin \vartheta_n t_n, \quad v_n = -\sin \vartheta_n s_n + \cos \vartheta_n t_n.$$

Straightforward calculations show that $\vartheta_n = \varphi_n - \varphi$ gives (7) with $\|u_n\| = \|v_n\| = 1 \forall n$. Also

$$\|Au_n - mu_n\|^2 \leq \cos^2 \vartheta_n \|As_n - ms_n\|^2 + \sin^2 \vartheta_n (M - m)^2,$$

and, since $\|As_n - ms_n\| \rightarrow 0$, $\vartheta_n \rightarrow 0$ as $n \rightarrow \infty$, $\|Au_n - mu_n\| \rightarrow 0$ as $n \rightarrow \infty$. Similarly, $\|Av_n - Mv_n\| \rightarrow 0$ as $n \rightarrow \infty$. \square

3. A theorem on Successive Transformations of a Probability Measure

THEOREM 2. *Let ν_0 be a probability measure on the family \mathcal{B} of Borel sets of $(0, \infty)$, with support $[m, M]$, so that*

$$m = \text{ess inf}(\nu_0) = \sup(\alpha; \nu_0\{x, x < \alpha\} = 0),$$

$$M = \text{ess sup}(\nu_0) = \inf(\alpha; \nu_0\{x, x > \alpha\} = 0).$$

Assume that $0 < m < M < \infty$. Consider the transformation $T: \nu_k \rightarrow \nu_{k+1}$ defined by

$$\nu_{k+1}(\mathcal{A}) = \int_{\mathcal{A}} \frac{(\lambda - \mu_1^k)^2}{D_k} \nu_k(d\lambda) \quad (8)$$

for any $\mathcal{A} \in \mathcal{B}$, where $\mu_1^k = \int \lambda \nu_k(d\lambda)$ and $D_k = \mu_2^k - (\mu_1^k)^2$, with $\mu_2^k = \int \lambda^2 \nu_k(d\lambda)$. Then, when $k \rightarrow \infty$,

$$\nu_{2k}(\mathcal{A}) \rightarrow p, \quad \nu_{2k+1}(\mathcal{A}) \rightarrow 1 - p \quad (9)$$

for all $\mathcal{A} = [m, x]$, $m < x < M$, for some p depending on ν_0 , $0 < p < 1$.

Proof. The proof is divided into five parts. In (i), we prove that the sequence of variances D_k is nondecreasing. In (ii), we construct sequences of intervals $\mathcal{L}_k = [m_k, m_k + \delta]$ and $\mathcal{R}_k = [M_k - \delta, M_k]$ in which the measure ν_k will tend to concentrate. In (iii) we prove that $\mathcal{R}_k \cap \mathcal{R}_{k+1} \neq \emptyset$ and in (iv) that the sequence M_k is nondecreasing. Finally, the limiting behaviour of ν_k is derived in (v).

(i) From the assumption $m < M$, $D_0 > 0$. Define the moments $\mu_l^k = \int \lambda^l \nu_k(d\lambda)$, $l = 0, 1, 2, \dots$, we have $D_{k+1} = \mu_2^{k+1} - (\mu_1^{k+1})^2$, with

$$\begin{aligned}\mu_1^{k+1} &= \int \lambda \frac{(\lambda - \mu_1^k)^2}{D_k} \nu_k(d\lambda) = \frac{(\mu_1^k)^3 - 2\mu_1^k \mu_2^k + \mu_3^k}{D_k}, \\ \mu_2^{k+1} &= \int \lambda^2 \frac{(\lambda - \mu_1^k)^2}{D_k} \nu_k(d\lambda) = \frac{(\mu_1^k)^2 \mu_2^k - 2\mu_1^k \mu_3^k + \mu_4^k}{D_k},\end{aligned}$$

so that

$$D_{k+1} - D_k = \frac{2\mu_1^k \mu_2^k \mu_3^k + \mu_2^k \mu_4^k - (\mu_1^k)^2 \mu_4^k - (\mu_3^k)^2 - (\mu_2^k)^3}{[\mu_2^k - (\mu_1^k)^2]^2} = \frac{\det M_3^k}{D_k^2},$$

where M_3^k is the moment matrix

$$M_3^k = \begin{pmatrix} \mu_0^k & \mu_1^k & \mu_2^k \\ \mu_1^k & \mu_2^k & \mu_3^k \\ \mu_2^k & \mu_3^k & \mu_4^k \end{pmatrix}.$$

Since M_3^k is nonnegative definite, $\det M_3^k \geq 0$ and therefore $D_{k+1} \geq D_k$. Since ν_k has bounded support, D_k is bounded from above, and therefore the sequence (D_k) monotonously converges to some limit D_* .

(ii) From part (i),

$$\det M_3^k = (D_{k+1} - D_k) D_k^2 \leq (D_{k+1} - D_k) D_*^2 \rightarrow 0 \quad \text{as } k \rightarrow \infty.$$

Therefore, given ϵ , $\exists K_\epsilon$ such that $\forall k > K_\epsilon$, $\det M_3^k < \epsilon$. Using Lemma 3 (see Appendix), for ϵ small enough, for any $k > K_\epsilon$ there exist two intervals $\mathcal{I}_1^k, \mathcal{I}_2^k$, with width $\delta = (M - m)\epsilon^{1/4}/(4D_0^{3/2})$ and such that $\nu_k(\mathcal{I}_1^k) + \nu_k(\mathcal{I}_2^k) \geq 1 - 4\sqrt{\epsilon}$, $\nu_k(\mathcal{I}_1^k) \geq D_k/(4M^2)$, $\nu_k(\mathcal{I}_2^k) \geq D_k/(4M^2)$. Also, $\max_{x \in \mathcal{I}_1^k, y \in \mathcal{I}_2^k} |x - y| \geq \sqrt{D_k/2}$. Without any loss of generality, assume that \mathcal{I}_1^k is the interval on the left. Define $\mathcal{L}(x) = [x, x + \delta]$, $\mathcal{R}(x) = [x - \delta, x]$,

$$\begin{aligned}\mathcal{X}_L^k &= \text{Arg max}_x \{ \nu_k[\mathcal{L}(x)], \mathcal{L}(x) \cap \mathcal{I}_1^k \neq \emptyset \}, \\ \mathcal{X}_R^k &= \text{Arg max}_x \{ \nu_k[\mathcal{R}(x)], \mathcal{R}(x) \cap \mathcal{I}_2^k \neq \emptyset \},\end{aligned}$$

and $m_k = \min \mathcal{X}_L^k$, $M_k = \max \mathcal{X}_R^k$, $\mathcal{L}_k = \mathcal{L}(m_k)$, $\mathcal{R}_k = \mathcal{R}(M_k)$; that is, M_k is the right endpoint of an interval \mathcal{R}_k , intersecting \mathcal{I}_2^k , with maximum measure, and similarly for m_k and \mathcal{L}_k . Note that $\nu_k(\mathcal{L}_k) + \nu_k(\mathcal{R}_k) \geq 1 - 4\sqrt{\epsilon}$, $\nu_k(\mathcal{L}_k) \geq D_k/(4M^2)$

and $v_k(\mathcal{R}_k) \geq D_k/(4M^2)$. The situation is the same for the two sequences of intervals (\mathcal{L}_k) and (\mathcal{R}_k) , and we concentrate on (\mathcal{R}_k) in the rest of the proof.

(iii) We show now that $\mathcal{R}_k \cap \mathcal{R}_{k+1} \neq \emptyset$. Again for ϵ small enough $\mu_1^k \notin \mathcal{R}_k$ and $\lambda - \mu_1^k \geq M_k - \delta - \mu_1^k$ on \mathcal{R}_k so that

$$\begin{aligned} v_{k+1}(\mathcal{R}_k) &= \int_{\mathcal{R}_k} \frac{(\lambda - \mu_1^k)^2}{D_k} v_k(d\lambda) \geq \frac{v_k(\mathcal{R}_k)}{D_k} (M_k - \delta - \mu_1^k)^2 \\ &\geq \frac{1}{4M^2} (M_k - \delta - \mu_1^k)^2. \end{aligned}$$

From Lemma 3,

$$\begin{aligned} M_k - \mu_1^k + \delta &\geq \frac{D_k^3}{2M(D_k^2 + D_k M^2 + M^4)} \geq \frac{D_0^3}{2M(D_*^2 + D_* M^2 + M^4)} = C, \end{aligned} \quad (10)$$

choosing ϵ such that $\delta < C/4$ gives $v_{k+1}(\mathcal{R}_k) > C/(8M^2)$. Choosing now ϵ such that $4\sqrt{\epsilon} < C/(8M^2)$ we obtain $\mathcal{R}_k \cap \mathcal{R}_{k+1} \neq \emptyset$ for any $k > K_\epsilon$.

(iv) We prove now that the sequence (M_k) is not decreasing starting at some K_ϵ for ϵ small enough. Take $k > K_\epsilon$ and assume that $M_{k+1} = M_k - \beta$, $\beta > 0$. Then note that $\beta < \delta$ since $\mathcal{R}_k \cap \mathcal{R}_{k+1} \neq \emptyset$ by (iii) above. Consider the difference $v_{k+1}(\mathcal{R}_k) - v_{k+1}(\mathcal{R}_{k+1}) = v_{k+1}([M_k - \beta, M_k]) - v_{k+1}([M_k - \delta - \beta, M_k - \delta])$. Assume first that $v_{k+1}([M_k - \delta - \beta, M_k - \delta]) = 0$, then $v_{k+1}(\mathcal{R}_k) > v_{k+1}(\mathcal{R}_{k+1})$, which is impossible by construction. We can thus consider the following ratio

$$\begin{aligned} &\frac{v_{k+1}([M_k - \beta, M_k])}{v_{k+1}([M_k - \delta - \beta, M_k - \delta])} \\ &= \frac{\int_{M_k - \beta}^{M_k} (\lambda - \mu_1^k)^2 v_k(d\lambda)}{\int_{M_k - \delta - \beta}^{M_k - \delta} (\lambda - \mu_1^k)^2 v_k(d\lambda)} \\ &\geq \frac{(M_k - \beta - \mu_1^k)^2}{(M_k - \delta - \mu_1^k)^2} \frac{v_k([M_k - \beta, M_k])}{v_k([M_k - \delta - \beta, M_k - \delta])} \geq 1. \end{aligned}$$

The last inequality follows from $\beta < \delta$, $M_k - \mu_1^k \geq 3\delta$, see (10) with $C > 4\delta$, and from

$$v_k([M_k - \beta, M_k]) \geq v_k([M_k - \delta - \beta, M_k - \delta])$$

which is a consequence of the definition of M_k noting that $v_k([M_k - \delta, M_k]) \geq v_k([M_k - \delta - \beta, M_k - \beta])$ and removing the interval $(M_k - \delta, M_k - \beta)$. Therefore, $\beta > 0$ leads to $v_{k+1}(\mathcal{R}_k) > v_{k+1}(\mathcal{R}_{k+1})$, which is impossible. We thus obtain $M_{k+1} \geq M_k$ for $k > K_\epsilon$.

(v) Since the sequence (M_k) is nondecreasing and bounded from above (by M), it has a limit $M_* \leq M$. The same is true for m_k , and $m_k \rightarrow m_*$ as $k \rightarrow \infty$. We have thus proved that for any δ small enough and any k larger than some K_δ ,

$$v_k([M_* - \delta, M_*]) + v_k([m_*, m_* + \delta]) \geq 1 - 64D_0^3\delta^2/(M - m)^2.$$

Assume that $M_* < M$. This would imply $\nu_k([M - \delta, M]) \rightarrow 0$ as $k \rightarrow \infty$ for $\delta < M - M_*$. On the other hand,

$$\frac{\nu_{k+1}([M - \delta, M])}{\nu_{k+1}([M_* - \delta, M_*])} > \frac{\nu_k([M - \delta, M])}{\nu_k([M_* - \delta, M_*])},$$

which leads to a contradiction since $\nu_k([M - \delta, M])/\nu_k([M_* - \delta, M_*])$ is then increasing and $\nu_k([M_* - \delta, M_*])$ is bounded from below. Therefore, $M_* = M$, and similarly $m_* = m$, with, for δ small enough and any k larger than some K_δ , $\nu_k([m + \delta, M - \delta]) < 64D_0^3\delta^2/(M - m)^2$. Finally, from Helly's theorem, see [8], p. 319, from the sequence (ν_k) we can extract a subsequence (ν_{k_i}) that is weakly convergent, and from the result above the associated limit has necessarily the form ν_p^* , where ν_p^* is the discrete measure concentrated on the two points m, M , with $\nu_p^*(m) = p$, $\nu_p^*(M) = 1 - p$. Since D_{k_i} converges to some D_* , ν_p^* is such that the associated variance is equal to D_* , which only leaves two possibilities for p (and $1 - p$):

$$1 - p, \quad p = \frac{1}{2} \pm \sqrt{\frac{1}{4} - \frac{D_*}{(M - m)^2}}.$$

Applying the transformation T , we get $\nu_{k_i+1} = T(\nu_{k_i}) \rightarrow T(\nu_p^*) = \nu_{1-p}^*$. \square

4. Stability of Attractors

In this section we consider the range of possible values for p in the attraction Theorem 1. We shall use the following definition of stability, see [3], p. 444, [5], p. 7.

DEFINITION 1. A fixed point ν^* for a mapping $T(\cdot)$ on probability measures on a family \mathcal{B} of Borel sets will be called stable if $\forall \epsilon > 0, \exists \alpha > 0$ such that for any ν_0 for which $d(\nu_0, \nu^*) < \alpha$, $d(T^n(\nu_0), \nu^*) < \epsilon$ for all $n > 0$. A fixed point ν^* is unstable if it is not stable.

We shall treat separately the stability and instability cases. For instability, we consider bounded quadratic operators in a Hilbert space and use the distance $d(\nu, \nu')$ given by the Lévy–Prokhorov metric, see [8], p. 349. In our case (measures supported on $[m, M]$), $d(\nu, \nu')$ becomes the Lévy distance between the distribution functions F, F' associated with ν, ν' , which we denote

$$L(F, F') = \inf\{\epsilon : F'(x - \epsilon) - \epsilon \leq F(x) \leq F'(x + \epsilon) + \epsilon, \quad \forall x\}.$$

In the case where one of the two measures is the discrete measure ν_p^* concentrated on m, M , with $\nu_p^*(m) = p$, $\nu_p^*(M) = 1 - p$, we get

$$d(\nu, \nu_p^*) = L(F, F_p^*) = \inf\{\epsilon : F(x) \leq p + \epsilon \text{ for } x < M - \epsilon \\ \text{and } p - \epsilon \leq F(x) \text{ for } m + \epsilon \leq x\},$$

with F_p^* the distribution function associated with v_p^* . We then have the following property.

THEOREM 3. *Consider the situation of Theorem 2, with v_0 any probability measure supported on some closed subset \mathcal{S}_A of $[m, M]$ and $\text{ess inf}(v_0) = m$, $\text{ess sup}(v_0) = M$.*

- (i) *The measure v_p^* is a fixed point for the mapping T^2 .*
- (ii) *Consider the set \mathcal{I}_u defined by*

$$\mathcal{I}_u = \left(0, \frac{1}{2} - s(\lambda^*)\right) \cup \left(\frac{1}{2} + s(\lambda^*), 1\right),$$

where

$$s(\lambda) = \frac{\sqrt{(M - \lambda)^2 + (\lambda - m)^2}}{2(M - m)}, \quad \lambda^* = \min_{\lambda \in \mathcal{S}_A} s(\lambda). \quad (11)$$

Any fixed point v_p^* with p in \mathcal{I}_u corresponds to an unstable fixed point for T^2 .

Proof. (i) It is straightforward to check that $T^2(v_p^*) = v_p^*$, $\forall p \in (0, 1)$.

(ii) We assume that \mathcal{S}_A is not reduced to $\{m, M\}$ (otherwise $\mathcal{I}_u = \emptyset$). One has $v_{k+2}(d\lambda) = H(v_k, \lambda)v_k(d\lambda)$, with

$$H(v_k, \lambda) = \frac{(\lambda - \mu_1^k)^2(\lambda - \mu_1^{k+1})^2}{D_k D_{k+1}},$$

see (8), with μ_1^k , D_k defined as in Theorem 2. For $v_k = v_p^*$, it gives

$$H(v_p^*, \lambda) = \frac{[M(1 - p) + mp - \lambda]^2 [Mp + m(1 - p) - \lambda]^2}{p^2(1 - p)^2(M - m)^4}.$$

One can then check that for any $p \in \mathcal{I}_u$, $\max_{\lambda \in \mathcal{S}_{v_0}} H(v_p^*, \lambda) = H(v_p^*, \lambda^*) > 1$, with $\lambda^* = \min_{\lambda \in \mathcal{S}_{v_0}} s(\lambda)$. Therefore, for any $p \in \mathcal{I}_u$, one can choose ϵ small enough, such that $d(v_k, v_p^*) < \epsilon$ implies $v_{k+2}([a, b]) > K_p v_k([a, b])$, for some $K_p > 1$ and some a, b such that $m + \epsilon < a < b < M - \epsilon$ and $[a, b] \cap \mathcal{S}_A \neq \emptyset$. For any $\alpha > 0$, $\alpha < 1 - p$, take an initial measure v_0 putting weight p at m , $1 - p - \alpha$ at M and α in the interval $[a, b]$. It satisfies $d(v_0, v_p^*) < \alpha$, and, for any m , either $d(v_{2m}, v_p^*) > \epsilon$ or $v_{2(m+1)}([a, b]) > K_p v_{2m}([a, b])$. The later case gives $v_{2m}([a, b]) > 2\epsilon$, and thus $d(v_{2m}, v_p^*) > \epsilon$, as soon as $m > \log(4\epsilon/\alpha)/\log(K_p)$, which shows that v_p^* is unstable. \square

One may notice that when v_0 is a discrete probability measure, the condition $H(v_p^*, \lambda) > 1$ corresponds to a condition on the eigenvalues of the Jacobian of the transformation T^2 , see [7]. Concerning stability, we restrict our attention to the finite-dimensional case $\mathcal{H} = \mathbb{R}^d$.

THEOREM 4. *Consider the situation of Theorem 2, with v_0 any discrete probability measure supported on $\{\lambda_1, \dots, \lambda_d\}$ with $0 < \lambda_1 < \lambda_2 \leq \dots \leq \lambda_{d-1} < \lambda_d$. Any fixed point v_p^* with p in the interval*

$$\mathcal{I}_S = \left(\frac{1}{2} - s(\lambda_{i^*}), \frac{1}{2} + s(\lambda_{i^*}) \right), \quad (12)$$

where $s(\lambda)$ is given by (11) and i^* is such that $|\lambda_{i^*} - (\lambda_1 + \lambda_d)/2|$ is minimum over all λ_i 's, $i = 2, \dots, d-1$, is stable for the mapping T^2 .

Proof. Consider the $(d-1)$ -dimensional canonical simplex $\mathcal{S}_{d-1} = \{w = (w_1, \dots, w_{d-1}) \mid w_i \geq 0, \sum_{i=1}^{d-1} w_i \leq 1\}$ and define $w_i^{(k)} = [z_k]_i^2$, $i = 1, \dots, d$. The transformation T of Theorem 2 is equivalently defined by (5), and T^2 defines an operator

$$\begin{aligned} \phi(\cdot): \mathcal{S}_{d-1} &\mapsto \mathcal{S}_{d-1} \quad \text{which maps} \\ (w_1^{(k)}, \dots, w_{d-1}^{(k)}) &\text{ to } (w_1^{(k+2)}, \dots, w_{d-1}^{(k+2)}). \end{aligned}$$

Studying the stability properties of T^2 is thus equivalent to studying those of ϕ . The transformation ϕ is defined by

$$\begin{aligned} w_i^{(k+2)} &= w_i^{(k+1)} \frac{(\sum_{j=1}^d \lambda_j w_j^{(k+1)} - \lambda_i)^2}{D_{k+1}} \\ &= w_i^{(k)} \frac{(\mu_1^k - \lambda_i)^2}{D_k} \frac{1}{D_{k+1}} \left(\frac{(\mu_1^k)^3 - 2\mu_1^k \mu_2^k + \mu_3^k}{D_k} - \lambda_i \right)^2, \\ &\quad i = 1, \dots, d, \end{aligned} \quad (13)$$

with D_k, μ_i^k defined as in Theorem 2, which gives

$$\begin{aligned} D_k &= \mu_2^k - (\mu_1^k)^2, \\ D_{k+1} &= \frac{1}{D_k} [(\mu_1^k)^2 \mu_2^k - 2\mu_1^k \mu_3^k + \mu_4^k] - \frac{1}{D_k^2} [(\mu_1^k)^3 - 2\mu_1^k \mu_2^k + \mu_3^k]^2. \end{aligned}$$

Take $w^{(k)} = (w_1^{(k)}, \dots, w_{d-1}^{(k)})$ in \mathcal{S}_{d-1} such that $w_1^{(k)} \in \mathcal{I}_S$ and $|w_i^{(k)}| < \alpha$, $i = 2, \dots, d-1$. The two-step iteration (13) can be written as

$$w_i^{(k+2)} = w_i^{(k)} H(\lambda_i, \mu_1^k, \mu_2^k, \mu_3^k, \mu_4^k), \quad i = 1, \dots, d-1. \quad (14)$$

Define $\bar{w}^{(k)} = (w_1^{(k)}, 0, \dots, 0) \in \mathcal{S}_{d-1}$ and $\bar{\mu}_m^k = \lambda_1^m w_1^{(k)} + \lambda_d^m (1 - w_1^{(k)})$, $m = 1, \dots, 4$. Then, for $i = 2, \dots, d-1$,

$$w_i^{(k+2)} = w_i^{(k)} h(\lambda_i, w_1^{(k)}) [1 + O(\alpha)], \quad \alpha \rightarrow 0,$$

with $h(\lambda, w_1^{(k)}) = H(\lambda, \bar{\mu}_1^k, \bar{\mu}_2^k, \bar{\mu}_3^k, \bar{\mu}_4^k)$. Straightforward calculations show that $h[\lambda, 1/2 \pm s(\lambda)] = 1$, with $s(\lambda)$ given by (11), and $w_1^{(k)} \in \mathcal{I}_S$ implies $h(\lambda_i, w_1^{(k)}) < 1$, $i = 2, \dots, d-1$. Therefore, for α smaller than some α_0 , there exist $L^{(k)} < 1$ such that

$$w_i^{(k+2)} < L^{(k)} w_i^{(k)}, \quad i = 2, \dots, d-1.$$

Similarly, (14) gives

$$w_1^{(k+2)} = w_1^{(k)} H(\lambda_1, \bar{\mu}_1^k, \bar{\mu}_2^k, \bar{\mu}_3^k, \bar{\mu}_4^k) + \sum_{j=2}^{d-1} w_1^{(k)} w_j^{(k)} \frac{\partial H(\lambda_1, \mu'_1, \mu'_2, \mu'_3, \mu'_4)}{\partial w_j} \Big|_{\bar{w}^{(k)}},$$

with $\tilde{w}^{(k)} = (1 - \gamma)\bar{w}^{(k)} + \gamma w^{(k)}$ for some $\gamma \in [0, 1]$, and $\mu'_m = \sum_{i=1}^{d-1} \lambda_i^m w_i + \lambda_d^m (1 - \sum_{i=1}^{d-1} w_i)$, $m = 1, \dots, 4$. Now, $H(\lambda_1, \bar{\mu}_1^k, \bar{\mu}_2^k, \bar{\mu}_3^k, \bar{\mu}_4^k) = 1$ and direct calculation shows that

$$\begin{aligned} & \frac{\partial H(\lambda_1, \mu'_1, \mu'_2, \mu'_3, \mu'_4)}{\partial w_j} \Big|_{\bar{w}^{(k)}} \\ &= \frac{(\lambda_j - \lambda_1)(\lambda_d - \lambda_j)^2 [2\lambda_d(1 - w_1^{(k)}) + 2\lambda_1 w_1^{(k)} - \lambda_1 - \lambda_j]}{(w_1^{(k)})^2 (1 - w_1^{(k)})^2 (\lambda_d - \lambda_1)^4} \end{aligned}$$

is positive for $w_1^{(k)} = 1/2 - s(\lambda_j)$, negative for $w_1^{(k)} = 1/2 + s(\lambda_j)$ and its absolute value reaches its maximum at one of these two points. By continuity arguments, the same is true for $\partial H(\lambda_1, \mu'_1, \mu'_2, \mu'_3, \mu'_4)/\partial w_j|_{\tilde{w}^{(k)}}$. By choosing α small enough, we thus guarantee that $w_1^{(k+2)} \in \mathcal{I}_S$, with

$$\min\{w_1^{(k)}, 1/2 - s(\lambda_{i^*}) + \delta\} \leq w_1^{(k+2)} \leq \max\{w_1^{(k)}, 1/2 + s(\lambda_{i^*}) - \delta\}$$

for some $\delta > 0$. This ensures that there exist $L < 1$ such that $L^{(k+2)} < L^* = \max(L^{(k)}, L) < 1$. Repeating the same arguments, we thus obtain

$$w_i^{(k+2m)} < (L^*)^m w_i^{(k)}, \quad i = 2, \dots, d-1.$$

Moreover, $|w_1^{(k+2)} - w_1^{(k)}| < B \sum_{i=2}^{d-1} w_i^{(k)}$ for some constant B , and therefore

$$|w_1^{(k+2m)} - w_1^{(k)}| < B \sum_{i=2}^{d-1} \sum_{n=0}^{m-1} w_i^{(k+2n)}$$

so that $|w_1^{(k+2m)} - w_1^{(k)}| < B(d-2)\alpha L^*/(1 - L^*)$ for any m . \square

Note that $\forall \lambda_1, \dots, \lambda_d$ the stability interval \mathcal{I}_S contains the interval

$$\left] \frac{1}{2} - \frac{1}{2\sqrt{2}}, \frac{1}{2} + \frac{1}{2\sqrt{2}} \right[\approx]0.14645, 0.85355[.$$

Numerical simulations for $d = 3$ show that for any initial density of x_1 in \mathbb{R}^d associated with a density of $w^{(1)}$ on S_{d-1} reasonably spread, the density of the values of p corresponding to stable attractors ν_p^* can be approximated by

$$\varphi(p) = C \log \min\{1, (J_\phi)_{22}\} = \begin{cases} C \log(J_\phi)_{22} & \text{if } p \in \mathcal{I}_S, \\ 0 & \text{otherwise,} \end{cases}$$

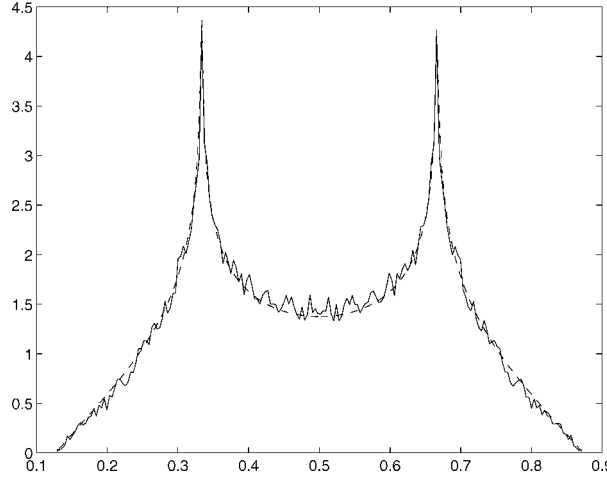


Figure 1. Empirical density of attractors (full line) and $\varphi(p)$ (dashed line) ($d = 3, \lambda_1 = 1, \lambda_2 = 4, \lambda_3 = 10$).

where C is a normalisation constant and J_ϕ is the Jacobian matrix of $\phi(\cdot)$ at points $w(p) = (p, 0, \dots, 0)$. It can be computed analytically, and is given by:

$$(J_\phi)_{ij} = \frac{\partial \phi_i(w)}{\partial w_j} \Big|_{w=w(p)} \begin{cases} 1 & \text{if } i = j = 1, \\ \frac{(\lambda_j - \lambda_1)(\lambda_d - \lambda_j)^2 [2\lambda_d(1-p) + 2\lambda_1 p - \lambda_1 - \lambda_j]}{p(1-p)^2(\lambda_d - \lambda_1)^4} & \text{if } j > 1 \text{ and } i = 1, \\ \frac{[\lambda_d(1-p) + \lambda_1 p - \lambda_j]^2 [\lambda_d p + \lambda_1(1-p) - \lambda_j]^2}{p^2(1-p)^2(\lambda_d - \lambda_1)^4} & \text{if } i = j, \\ 0 & \text{otherwise.} \end{cases}$$

Figure 1 shows the empirical density of attractors (full line) together with $\varphi(p)$ (dashed line) in the case $\lambda_1 = 1, \lambda_2 = 4, \lambda_3 = 10$. The support of this density coincides with the stability interval \mathcal{I}_S given by (12). When $d > 3$, the density of attractors depends on the initial density of x_1 .

5. Rate of Convergence

Define the convergence rate at iteration k by

$$r_k = \frac{f(x_{k+1}) - f(x^*)}{f(x_k) - f(x^*)}, \quad (15)$$

where $f(x)$ is given by (1). One can rewrite r_k as

$$r_k = \frac{(A^{-1}g_{k+1}, g_{k+1})}{(A^{-1}g_k, g_k)}$$

and using the updating rule for g_k we get

$$r_k = 1 - \frac{(g_k, g_k)^2}{(Ag_k, g_k)(A^{-1}g_k, g_k)} = 1 - \frac{1}{\mu_1^k \mu_{-1}^k}$$

with μ_m^k defined as in Theorem 2. Define the asymptotic rate as

$$R = R(x_1, x^*) = \lim_{k \rightarrow \infty} \left(\prod_{j=1}^k r_j \right)^{1/k}. \quad (16)$$

Generally, R depends on the initial point x_1 and the optimal point x^* . Theorem 1 implies that for any fixed x^* and almost all x_1 the asymptotic rate only depends on the value of p and is given by $R(p) = \{[f(x_{k+2}) - f(x^*)]/[f(x_k) - f(x^*)]\}^{1/2}$, where x_k is associated with z_k given by (7). This gives

$$R(p) = \frac{p(1-p)(\rho-1)^2}{[p + \rho(1-p)][(1-p) + \rho p]},$$

with $\rho = M/m$. When $\mathcal{H} = \mathbb{R}^d$, $\rho = \lambda_d/\lambda_1$ is the condition number of the matrix A . The function $R(p)$ is symmetric with respect to $1/2$ and monotonously increasing from 0 to $1/2$. The worst asymptotic rate is thus obtained at $p = 1/2$:

$$R_{\max} = \left(\frac{\rho-1}{\rho+1} \right)^2. \quad (17)$$

Note that from the Kantorovich inequality, see [4] and [6], p. 151,

$$\mu_1^k \mu_{-1}^k \leq (1 + \rho)^2 / (4\rho),$$

and therefore $\forall x_k, r_k \leq R_{\max}$. In the finite-dimensional case, the worst rate is thus achieved only when $[x_k]_1 = \pm \rho[x_k]_d, [x_k]_2 = \dots = [x_k]_{d-1} = 0$.

Consider now another convergence rate, defined by

$$R' = \lim_{k \rightarrow \infty} \left(\prod_{j=1}^k r'_j \right)^{1/k},$$

where $r'_k = ([x_{k+1} - x^*], [x_{k+1} - x^*]) / ([x_k - x^*], [x_k - x^*])$. In a way similar to the derivation of r_k , by rewriting r'_k in terms of z_k , one gets

$$r'_k = 1 - \frac{2\mu_{-1}^k}{\mu_1^k \mu_{-2}^k} + \frac{1}{(\mu_1^k)^2 \mu_{-2}^k}.$$

One can easily check that for almost all x_1 the asymptotic rate R' is equal to $R(p)$, where p defines the attractor. Similarly, we can define R'' from $r''_k = (g_{k+1}, g_{k+1}) / (g_k, g_k)$, which gives $r''_k = \mu_2^k / (\mu_1^k)^2 - 1$ and again $R'' = R(p)$ for almost all x_1 .

Several definitions of the convergence rate give asymptotically the same value $R(p)$. According to the results of Section 4, only values of p in \mathcal{I}_S given by (12) may correspond to stable attractors. The range of possible values of $R(p)$ is thus $[R_{\min}, R_{\max}]$, where R_{\max} , given by (17), is obtained for $p = 1/2$ and

$$R_{\min} \leq R_{\min}^* = R(1/2 + 1/[2\sqrt{2}]) = \frac{(\rho-1)^2}{(\rho+1)^2 + 4\rho}.$$

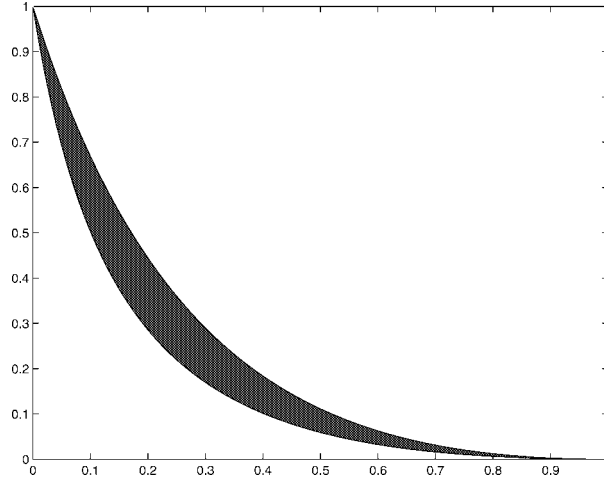


Figure 2. Range $[R_{\min}^*, R_{\max}]$ of possible values for the asymptotic rate $R(p)$ as a function of ρ .

Figure 2 presents the range $[R_{\min}^*, R_{\max}]$ as a function of ρ , the upper curve corresponding to R_{\max} and the lower to R_{\min}^* . The maximum value of the range is $3 - 2\sqrt{2} \simeq 0.1716$, obtained at $\rho = 1 + 2\sqrt{2} - 2\sqrt{2 + \sqrt{2}} \simeq 0.1329$. It confirms the experimental observation that the rate of convergence of the gradient algorithm is generally close to its worst value R_{\max} .

6. Steepest Descent with Relaxation

The introduction of a relaxation coefficient γ , with $0 < \gamma < 1$, in the steepest-descent algorithm totally changes its behaviour. The algorithm (2) then becomes

$$x_{k+1} = x_k - \gamma \frac{(g_k, g_k)}{(Ag_k, g_k)} g_k.$$

We restrict our attention to the finite dimensional case $\mathcal{H} = \mathbb{R}^d$. For fixed A , depending on the value of γ , the renormalized process either converges to periodic orbits (the same for almost all starting points) or exhibits a chaotic behaviour. Figures 3 presents the classical period-doubling phenomenon in the case $d = 2$ when $\lambda_1 = 1$ and $\lambda_2 = 10$. Figures 4 gives the asymptotic rate (16) as a function of γ in the same situation.

We get now instead of (15): $r_k(\gamma) = 1 - \gamma(2 - \gamma)/(\mu_1\mu_{-1})$. Note that from the Kantorovich inequality the worst value of the rate is

$$1 - \gamma(2 - \gamma) \frac{4\rho}{(1 + \rho)^2} > R_{\max},$$

if $\gamma < 1$. However, numerical results show that for γ large enough the asymptotic rate is significantly better than R_{\max} (see, for instance, Figure 4). A detailed analysis of the two-dimensional case gives the following results.

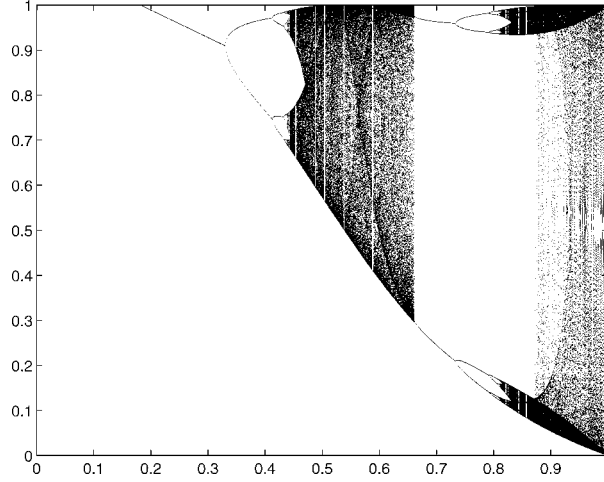


Figure 3. Attractors for z_1 as a function of γ ($d = 2, \lambda_1 = 1, \lambda_2 = 10$).

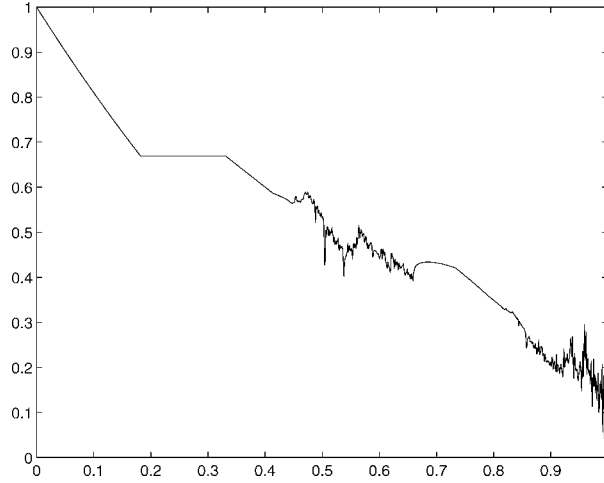


Figure 4. Asymptotic rate (16) as a function of γ ($d = 2, \lambda_1 = 1, \lambda_2 = 10$).
 $R_{\max} = (9/11)^2 \simeq 0.6694$.

- (i) If $0 < \gamma \leq 2/(\rho + 1)$, the process converges to the fixed point $p = 1$ and $R = R(\gamma) = 1 - 2\rho\gamma/(\rho + 1)$.
- (ii) If $2/(\rho + 1) < \gamma \leq 4\rho/(\rho + 1)^2$, the process converges to the fixed point

$$p = \frac{2\rho - \gamma(\rho + 1)}{2(\rho - 1)} \quad \text{and} \quad R(\gamma) = R_{\max}.$$

- (iii) If $4\rho/(\rho + 1)^2 < \gamma \leq 2(\sqrt{2} + 1)\rho/(\rho + 1)^2$, the process converges to the 2-point cycle (p_1, p_2) , with

$$p_{1,2} = \frac{2\rho - \gamma(\rho + 1) \pm \sqrt{\gamma(\gamma\rho^2 - 4\rho + 2\gamma\rho + \gamma)}}{2(\rho - 1)},$$

and $R(\gamma) = 1 - \gamma$.

- (iv) For larger values of γ one observes a classical period-doubling phenomenon, see Figure 3.
- (v) If $\rho > 3 + 2\sqrt{2} \approx 5.828427$, the process converges again to a 2-point cycle for values of γ larger than $\gamma_\rho = 8\rho/(\rho + 1)^2$, see Figure 3. For the limiting case $\gamma = \gamma_\rho$, the cycle is given by (p'_1, p'_2) , with

$$p'_{1,2} = \frac{\rho(\rho^2 - 2\rho + 5 \pm 2\sqrt{(\rho^2 - 2\rho + 5)(5\rho^2 - 2\rho + 1)})}{(\rho - 1)(\rho + 1)^3},$$

and the associated asymptotic rate is $R(\gamma_\rho) = (\rho^2 - 6\rho + 1)/(\rho^2 - 1)$.

In higher dimensions, repeated numerical trials show that the process typically no longer converges to the two-dimensional plane spanned by (e_1, e_d) .

Appendix

LEMMA 1. *Let v be any probability distribution on $[m, M]$, $0 < m \leq M < \infty$. Assume that $\exists \mathfrak{I}, |\mathfrak{I}| \leq \alpha$ and $v(\mathfrak{I}) \geq 1 - \epsilon$. Then, $\text{Var}(v) \leq \alpha^2 + 2\epsilon M^2$.*

Proof. Define $\mu_1 = \int \lambda v(d\lambda)$, $\mu_{\mathfrak{I}} = \int_{\mathfrak{I}} \lambda v(d\lambda)$. Then $\mu_1 = \mu_{\mathfrak{I}} + \int_{[m, M] \setminus \mathfrak{I}} \lambda v(d\lambda)$. Therefore, $\mu_{\mathfrak{I}} \leq \mu_1 \leq \mu_{\mathfrak{I}} + \epsilon M$. We get

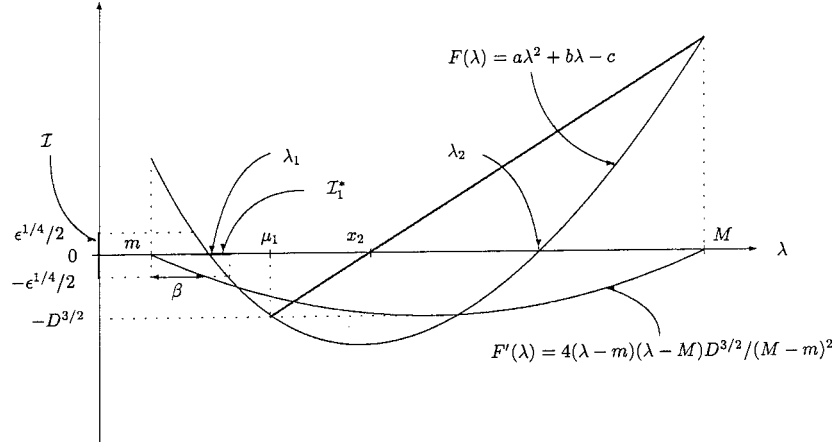
$$\begin{aligned} \text{Var}(v) &= \int (\lambda - \mu_1)^2 v(d\lambda) \\ &\leq \int_{\mathfrak{I}} (\lambda - \mu_1)^2 v(d\lambda) + (M - m)^2 \epsilon \\ &= \int_{\mathfrak{I}} (\lambda - \mu_{\mathfrak{I}})^2 v(d\lambda) + (1 - \epsilon)(\mu_1 - \mu_{\mathfrak{I}})^2 + (M - m)^2 \epsilon \\ &\leq \alpha^2(1 - \epsilon) + (1 - \epsilon)\epsilon^2 M^2 + (M - m)^2 \epsilon \\ &\leq \alpha^2 + 2\epsilon M^2. \end{aligned} \quad \square$$

LEMMA 2. *Let v be any probability distribution on $[m, M]$, $0 < m \leq M < \infty$. Assume that $\text{Var}(v) \leq \epsilon$. Then, there exist an interval \mathfrak{I} such that $|\mathfrak{I}| \leq \epsilon^{1/4}$ and $v(\mathfrak{I}) \geq 1 - 4\sqrt{\epsilon}$.*

Proof. Take $\mathfrak{I} = [\mu_1 - \epsilon^{1/4}/2, \mu_1 + \epsilon^{1/4}/2]$, $\mu_1 = \int \lambda v(d\lambda)$, and apply Chebyshev inequality. \square

LEMMA 3. *Let v be any distribution on $[m, M]$, $0 < m \leq M < \infty$. Define $\mu_n = \int \lambda^n v(d\lambda)$ and*

$$M_3 = \begin{pmatrix} \mu_0 & \mu_1 & \mu_2 \\ \mu_1 & \mu_2 & \mu_3 \\ \mu_2 & \mu_3 & \mu_4 \end{pmatrix}.$$

Figure 5. Construction of the intervals $\mathcal{J}_1, \mathcal{J}_2$.

Assume that $D = \text{Var}(v) = \mu_2 - \mu_1^2 > 0$ and $\det M_3 < \epsilon$. Then, there exist two intervals \mathcal{J}_1 and \mathcal{J}_2 such that

$$(i) \quad |\mathcal{J}_i| \leq \frac{(M - m)\epsilon^{1/4}}{4D^{3/2}}, \quad i = 1, 2, \quad v(\mathcal{J}_1) + v(\mathcal{J}_2) \geq 1 - 4\sqrt{\epsilon},$$

$$(ii) \quad \max_{x \in \mathcal{J}_i} |x - \mu_1| > \frac{D^3}{2M(D^2 + DM^2 + M^4)}, \quad i = 1, 2, \quad (18)$$

$$(iii) \quad \text{for } \epsilon < \epsilon_* = \frac{8D^4}{128D^3M^2 + (M - m)^2}, \quad v(\mathcal{J}_i) \geq \frac{D}{4M^2}, \quad i = 1, 2,$$

$$\text{and } \max_{x \in \mathcal{J}_1, y \in \mathcal{J}_2} |x - y| > \sqrt{D/2}. \quad (19)$$

Proof. (i) Define $a = \sqrt{D}$, $b = (\mu_1\mu_2 - \mu_3)/\sqrt{D}$, $c = a\mu_2 + b\mu_1 = (\mu_2^2 - \mu_1\mu_3)/\sqrt{D}$. Note that $a > 0$, $b < 0$ and $c < 0$. Define $\eta = F(\lambda) = a\lambda^2 + b\lambda - c$, with λ having the distribution v , so that $E\{\eta\} = 0$ and $\det M_3 = \text{Var}(\eta)$. From Lemma 2, the interval $\mathcal{I} = [-\epsilon^{1/4}/2, \epsilon^{1/4}/2]$ is such that $\text{Prob}\{\eta \in \mathcal{I}\} \geq 1 - 4\sqrt{\epsilon}$. From the mean-value theorem, there exist $\lambda_1 < \lambda_2$ such that $\lambda_i \in [m, M]$ and $a\lambda_i^2 + b\lambda_i - c = 0$, $i = 1, 2$. Direct calculation gives $F(\mu_1) = a\mu_1^2 + b\mu_1 - c = -D^{3/2}$. Take $\mathcal{J}_i = [\lambda_i - \beta, \lambda_i + \beta]$, $\beta = (M - m)\epsilon^{1/4}/(8D^{3/2})$, $i = 1, 2$. The two intervals $\mathcal{J}_1^*, \mathcal{J}_2^*$ that are mapped to \mathcal{I} by F : $\lambda \mapsto a\lambda^2 + b\lambda - c$ are such that $\mathcal{J}_i^* \subset \mathcal{J}_i$, $i = 1, 2$, and therefore $v(\mathcal{J}_1) + v(\mathcal{J}_2) \geq 1 - 4\sqrt{\epsilon}$. The fact that $\mathcal{J}_i^* \subset \mathcal{J}_i$ can be seen from the following: $m < \lambda_1 < \lambda_2 < M$, f is a quadratic function such that $F(\lambda_i) = 0$, $i = 1, 2$, with minimum less than $F(\mu_1) = -D^{3/2}$, see Figure 5 where $|\mathcal{J}_1^*| < 2\beta' < 2\beta$, with β' such that $F'(m + \beta') = -\epsilon^{1/4}/2$, $F'(\lambda) = 4(\lambda - m)(\lambda - M)D^{3/2}/(M - m)^2$.

(ii) Consider the interval \mathcal{J}_2 . We have $\max_{x \in \mathcal{J}_2} |x - \mu_1| > |\lambda_2 - \mu_1| > |x_2 - \mu_1|$, with x_2 such that the segment joining the points $\{\mu_1, F(\mu_1)\}$ and $\{M, F(M)\}$ crosses the horizontal axis $F = 0$ at x_2 , see Figure 5. This value

satisfies $x_2 - \mu_1 = D^{3/2}(M - \mu_1)/[F(M) + D^{3/2}]$. Since

$$D = \mu_2 - \mu_1^2 < M^2 - \mu_1^2 = (M - \mu_1)(M + \mu_1) < 2M(M - \mu_1),$$

$(M - \mu_1) > D/(2M)$. Also, $F(M) = \sqrt{D}M^2 + bM - c$ with $b < 0$ and $c > -M^4/\sqrt{D}$. Therefore, $F(M) < \sqrt{D}M^2 + M^4/\sqrt{D}$ and

$$x_2 - \mu_1 > \frac{D^3}{2M(D^2 + DM^2 + M^4)}.$$

Similarly, in the case of \mathfrak{l}_1 we also have $F(m) < \sqrt{D}M^2 + M^4/\sqrt{D}$, and by symmetry, using the transformation $\lambda \rightarrow \lambda' = M + m - \lambda$, we get $(\mu_1 - m) > D/(2M)$, which gives (18).

(iii) Assume that $v(\mathfrak{l}_2) < \gamma$, part (i) implies $v(\mathfrak{l}_1) > 1 - 4\sqrt{\epsilon} - \gamma$, and from Lemma 1

$$D \leq \frac{(M - m)^2 \sqrt{\epsilon}}{16D^3} + 2(4\sqrt{\epsilon} + \gamma)M^2,$$

which gives

$$\gamma \geq \frac{D}{2M^2} - \sqrt{\epsilon} \left[\frac{(M - m)^2}{32D^3M^2} + 4 \right],$$

and thus $\gamma \geq D/(4M^2)$ for $\epsilon < \epsilon_*$, see (19).

Define now $\Delta = \max_{x \in \mathfrak{l}_1, y \in \mathfrak{l}_2} |x - y|$. Lemma 1 gives $D \leq \Delta^2 + 8\sqrt{\epsilon}M^2$, which implies $\Delta^2 \geq D^2 - 8M^2\sqrt{\epsilon}$. Since $\epsilon < \epsilon_*$ implies $\epsilon < [D/(16M^2)]^2$, we get $\Delta^2 > D/2$. \square

References

1. Akaike, H.: On a successive transformation of probability distribution and its application to the analysis of the optimum gradient method, *Ann. Inst. Statist. Math. Tokyo* **11** (1959), 1–16.
2. Forsythe, G. E.: On the asymptotic directions of the s -dimensional optimum gradient method, *Numerische Math.* **11** (1968), 57–76.
3. Hale, J. and Koçak, H.: *Dynamics and Bifurcations*, Springer-Verlag, Heidelberg, 1991.
4. Kantorovich, L. V. and Akilov, G. P.: *Functional Analysis*, 2nd edn, Pergamon Press, London, 1982.
5. LaSalle, J. P.: *The Stability of Dynamical Systems*, SIAM, Philadelphia, 1976.
6. Luenberger, D. G.: *Introduction to Linear and Nonlinear Programming*, Addison-Wesley, Reading, Mass., 1973.
7. Pronzato, L., Wynn, H. P. and Zhigljavsky, A. A.: *Dynamical Search*, Chapman & Hall/CRC, Boca Raton, 2000.
8. Shiryaev, A. N.: *Probability*, Springer-Verlag, Berlin, 1996.