

A linear convergence proof of the isotropic ES via a continuous time approximation

Youhei Akimoto Anne Auger Nikolaus Hansen

TAO - INRIA Saclay

Oct. 19, 2012

Motivation

- ▶ Adaptive Evolution Strategies sample from a Gaussian distribution and adapt the mean vector and the covariance matrix
- ▶ Adaptive Evolution Strategies including CMA-ES are successfully applied in practice and there are lots of empirical evidence of linear convergence on a wide class of functions.
- ▶ However, their linear convergence to a local optimum is so far only proven for simple algorithms compared to CMA-ES [Auger 05, Jagerskupper 06, 07].

Motivation

- ▶ In this work, we seek a new methodology to show linear convergence.
 - ▶ *theory of stochastic approximation (ODE method)*
- ▶ An isotropic variant of adaptive-ES derived from IGO is studied.
- ▶ We derive
 - ▶ global linear convergence in expectation on monotonic convex-quadratic-composite functions
 - ▶ local linear convergence w.p. $1 - \varepsilon$ on monotonic \mathcal{C}^2 -composite functions

Information-Geometric Optimization (IGO) Algorithm

Choose a probabilistic model $\{P_\theta, \theta \in \Theta\}$ on \mathbb{X} and initialize θ_0 .

1. Draw samples $x_i, i = \llbracket 1, \lambda \rrbracket$, independently from P_{θ_n} .
2. Evaluate $f(x_i)$ and compute ranking
 $\text{rk}(x_i) = \#\{j \in \llbracket 1, \lambda \rrbracket : f(x_j) \leq f(x_i)\}$.
3. Update parameters as

$$\theta_{n+1} = \theta_n + \eta \sum_{i=1}^{\lambda} \underbrace{w_{\text{rk}(x_i)}}_{\text{nonincreasing weights}} \overbrace{\tilde{\nabla}_\theta \ln P_{\theta_n}(x_i)}^{\text{natural gradient of log-likelihood}} .$$

4. Go back to step 1 unless a termination criterion is satisfied.

\implies **PBIL** if P_θ is Bernoulli on $\mathbb{X} = \{0, 1\}^d$
 \implies **Pure rank- μ CMA-ES** if P_θ is Gaussian $\mathcal{N}(m, C)$ on $\mathbb{X} = \mathbb{R}^d$

Isotropic ES

= IGO algorithm with $P_\theta = \mathcal{N}(m, vI_d)$ on $\mathbb{X} = \mathbb{R}^d$.

Initialize $\theta_0 = (m_0, v_0)$, $m_0 \in \mathbb{R}^d$ and $v_0 > 0$.

1. Draw samples x_i , $i = \llbracket 1, \lambda \rrbracket$, independently from $\mathcal{N}(m_n, v_n I_d)$.
2. Evaluate $f(x_i)$ and compute ranking
 $\text{rk}(x_i) = \#\{j \in \llbracket 1, \lambda \rrbracket : f(x_j) \leq f(x_i)\}$.
3. Update parameters as

$$m_{n+1} = m_n + \eta_m \sum_{i=1}^{\lambda} w_{\text{rk}(x_i)} (x_i - m_n)$$

$$v_{n+1} = v_n + \eta_v \sum_{i=1}^{\lambda} w_{\text{rk}(x_i)} \left(\frac{\|x_i - m_n\|^2}{d} - v_n \right).$$

4. Go back to step 1 unless a termination criterion is satisfied.

Invariance of Isotropic ES

Invariance under monotone transformation of f

$f(x) \mapsto g(f(x))$, $g : \mathbb{R} \rightarrow \mathbb{R}$ is a strictly increasing function.

Invariance under translation of \mathbb{X}

$x \mapsto a + x$, $a \in \mathbb{R}^d$.

Invariance under rotation of \mathbb{X}

$x \mapsto Rx$, $R \in \mathbb{R}^{d \times d}$ is an orthogonal matrix.

Monotonic Convex Quadratic Composite Functions

Objective function $f(x) = g((x - \mathbf{x}^*)^T A(x - \mathbf{x}^*))$

- ▶ $g : \mathbb{R} \rightarrow \mathbb{R}$ is a strictly increasing function.
- ▶ $A \in \mathbb{R}^{d \times d}$ is a positive definite symmetric matrix.
- ▶ $\mathbf{x}^* \in \mathbb{R}^d$ is the global minimum.

W.l.o.g. (thanks to the invariance properties), we assume

- ▶ g is the identity function
- ▶ $A = \text{diag}(A_1, \dots, A_d)$
- ▶ $\mathbf{x}^* = 0$.

Monotonic \mathcal{C}^2 -Composite Functions

Objective function $f(x) = g(h(x))$

- ▶ $g : \mathbb{R} \rightarrow \mathbb{R}$ is a strictly increasing function.
- ▶ $h : \mathbb{R}^d \rightarrow \mathbb{R}$ is a \mathcal{C}^2 function.
- ▶ $\mathbf{x}^* \in \mathbb{R}^d$ is a local minimum of h .
- ▶ $A \in \mathbb{R}^{d \times d}$ is a positive definite symmetric Hessian of h at \mathbf{x}^* .

W.l.o.g. (thanks to the invariance properties), we assume

- ▶ g is the identity function
- ▶ $A = \text{diag}(A_1, \dots, A_d)$
- ▶ $\mathbf{x}^* = 0$.

Linear Convergence

aka log-linear convergence or geometric convergence

Definition (Linear convergence of x to \mathbf{x}^* in $(\mathbb{X}, \rho_{\mathbb{X}})$)

Let $\{x_n\}$ be a sequence in \mathbb{X} , $\mathbf{x}^* \in \mathbb{X}$ and $\rho_{\mathbb{X}}$ a metric on \mathbb{X} . If there exists $0 < \gamma_1 < 1$ such that

$$\limsup_{n \rightarrow \infty} \frac{\rho_{\mathbb{X}}(x_{n+1}, \mathbf{x}^*)}{\rho_{\mathbb{X}}(x_n, \mathbf{x}^*)} = \gamma_1$$

or if there exists $\gamma_2 < 0$ such that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \ln \rho_{\mathbb{X}}(x_n, \mathbf{x}^*) = \gamma_2 ,$$

$\{x_n\}$ is said to converge linearly (aka log-linearly or geometrically) to \mathbf{x}^* in $(\mathbb{X}, \rho_{\mathbb{X}})$.

Our Objective

- ▶ Expected distance of $x \sim P_{\theta_n}$ from the optimum \mathbf{x}^* given θ_n

$$\mathbb{E}_{x \sim P_{\theta_n}} [\|x - \mathbf{x}^*\|] \leq (\|m_n - \mathbf{x}^*\|^2 + d \cdot v_n)^{1/2}$$

- ▶ Optimum parameter $\theta^* \in \bar{\Theta}$

$$\theta^* = (\mathbf{x}^*, 0) \in \partial\Theta$$

- ▶ Metric on parameter space $\bar{\Theta}$

$$\rho(\theta, \theta') := \left[\|m - m'\|^2 + d(\sqrt{v} - \sqrt{v'})^2 \right]^{1/2}$$

$$\implies \rho(\theta, \theta^*) = (\|m_n - \mathbf{x}^*\|^2 + d \cdot v_n)^{1/2}$$

Linear convergence of θ_n to θ^* in $(\bar{\Theta}, \rho)$

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \ln \mathbb{E}_{x \sim P_{\theta_n}} [\|x - \mathbf{x}^*\|] \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \ln \rho(\theta_n, \theta^*) = \gamma < 0$$

w.p.1, in expectation, in probability, w.p. $1 - \varepsilon$, or etc.

Main Results

Under an assumption on w , (introduced later)

Global Linear Convergence on Monotonic Convex-Quadratic-Composite

For sufficiently small $\delta t > 0$, there exists $\gamma < 0$ such that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \ln \mathbb{E} [\rho(\theta_n, \theta^*)] = \gamma .$$

Local Linear Convergence on Monotonic \mathcal{C}^2 -Composite

For $\theta_0 \in U$ and any $0 < \varepsilon < 1$, there exists $\overline{\delta t} > 0$ such that for $0 < \forall \delta t \leq \overline{\delta t}$

$$\Pr \left[\limsup_{n \rightarrow \infty} \frac{1}{n} \ln \rho(\theta_n, \theta^*) \leq \gamma \right] \geq 1 - \varepsilon ,$$

where the speed of convergence $\gamma < 0$ depends on δt .

Notation

Let $\theta_n = [m_n^\top, v_n]^\top$ and

$$\hat{F}_n = \begin{bmatrix} \hat{F}_n^m \\ \hat{F}_n^v \end{bmatrix} = \frac{1}{\lambda} \sum_{i=1}^{\lambda} w_{\text{rk}(X_i)} \begin{bmatrix} x_i - m_n \\ \|x_i - m_n\|^2 / d - v_n \end{bmatrix} .$$

Then, the isotropic ES is $\theta_{n+1} = \theta_n + \delta t \hat{F}_n$ ($\eta_m = \eta_v = \delta t$).

Decompose

$$\hat{F}_n = F(\theta_n) + M_n ,$$

where

- ▶ $F(\theta) = \mathbb{E}[\hat{F}_n \mid \theta_n = \theta]$ is the *mean field*
- ▶ $M_n = \hat{F}_n - F(\theta_n)$ is the *martingale difference noise*

Then, the isotropic ES is $\theta_{n+1} = \theta_n + \delta t [F(\theta_n) + M_n]$.

Idea 1: Potential function

i.e., Drift function

Potential function

For some $\eta > 0$, define

$$\Psi_{\eta}(\theta) = \left(\eta \vee \frac{\|m\|}{\sqrt{v}} \right)^{1/2} \cdot \underbrace{(\|m\|^2 + d \cdot v)^{1/2}}_{=\rho(\theta, \theta^*)}$$

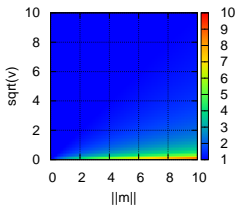
and show for some $N \in \mathbb{N}$ and $0 < \gamma < 1$

$$\mathbb{E}_n [\Psi_{\eta}(\theta_{N+n})] \leq \gamma \Psi_{\eta}(\theta_n) .$$

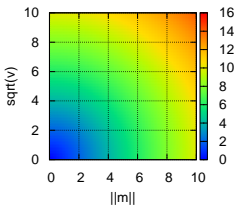
Note. Since $\rho(\theta, \theta^*) \leq \Psi_{\eta}(\theta_n) / \eta^{1/2}$,

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \ln \mathbb{E}_{\theta_n} [\rho(\theta, \theta^*)] \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \ln \mathbb{E}_{\theta_n} [\Psi_{\eta}(\theta_n)] .$$

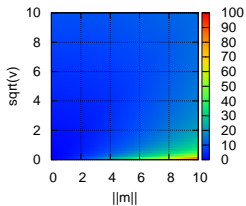
$$\eta \vee \frac{\|m\|}{\sqrt{v}}$$



$$\rho(\theta, \theta^*)$$



$$\Psi_\eta(\theta)$$



Idea 2: Stochastic Approximation

The isotropic ES as a noisy approximation of an ODE solution

Stochastic Approximation

Consider the isotropic ES

$$\theta_{n+1} = \theta_n + \delta t [F(\theta_n) + M_n]$$

as a noisy approximation of the solution to the ODE

$$\frac{d\theta}{dt} = F(\theta) .$$

In other words, θ_{N+n} is approximated by the solution $\varphi(N\delta t, \theta_n)$ of the ODE starting from θ_n .

Continuous time trajectory

Mean Field

The mean field $F(\theta)$ is expressed as

$$F(\theta) = \int w(P_\theta[y : f(y) \leq f(x)]) \begin{bmatrix} \alpha \tilde{\nabla}_{\theta_m} \ln p_\theta(x) \\ \tilde{\nabla}_{\theta_v} \ln p_\theta(x) \end{bmatrix} P_\theta(dx) ,$$

where $w : [0, 1] \rightarrow \mathbb{R}$ is a function defined by

$$w(p) = \sum_{i=1}^{\lambda} w_i \binom{\lambda-1}{i-1} p^{i-1} (1-p)^{\lambda-i} .$$

Global Existence and Uniqueness

The autonomous system $\frac{d\theta}{dt} = F(\theta)$ with initial value $\theta(0) = \theta_0$ has a unique solution on $[0, \infty)$ for each $\theta_0 \in \Theta$, i.e. there is only one solution $\varphi(\cdot, \theta_0) : t \rightarrow \theta(t)$ to the system for any $\theta_0 \in \Theta$.

Outline of the Proof

1. Bound the progress of $\varphi(N\delta t, \theta_n)$ over $\theta_n = \varphi(0, \theta_n)$.
2. Bound N -step cumulative error between θ_{n+N} and $\varphi(N\delta t, \theta_n)$.
3. Combine them and show $\mathbb{E}_n [\Psi_\eta(\theta_{N+n})] \leq \gamma \Psi_\eta(\theta_n)$.

Step 1: Study of Continuous-time Trajectory

Assumption on w

$$\int w(P_1[y : y \leq z])(z^2 - 1)P_1(dz) = \alpha > 0 ,$$

equivalently, $F^v(\theta)/v = \alpha/d$ for $f(x) = e^T x, \forall e \in \mathbb{R}^d \setminus \{0\}$.

Lemma 1: Exponential Convergence

There exist constants $\eta > 0$ and $\kappa_e > 0$ such that

$$\Psi_\eta(\varphi(t, \theta_0)) \leq \Psi_\eta(\theta_0)e^{-\kappa_e t}, \quad \forall t \geq 0, \forall \theta_0 \in \Theta.$$

For any $\theta_n \in \Theta$ and $t = N\delta t$, we have

$$\Psi_\eta(\varphi(N\delta t, \theta_n)) \leq \Psi_\eta(\theta_n)e^{-\kappa_e N\delta t} .$$

Proof of Lemma 1

Rewrite $\Psi_\eta(\theta) = \sqrt{\eta \nabla \zeta(\theta)} \sqrt{V(\theta)}$, where

- ▶ $\zeta(\theta) = \|m\| / \sqrt{v}$
- ▶ $V(\theta) = \|m\|^2 + d \cdot v = \rho(\theta, \theta^*)^2$

To show exponential convergence, we need to show

$$\frac{d\Psi_\eta(\theta)}{dt} = \nabla \Psi_\eta(\theta)^T F(\theta) \leq -\kappa_e \Psi_\eta(\theta) .$$

Proposition 1

There exists a constant $\kappa_V > 0$ such that

$$\nabla V(\theta)^T F(\theta) \leq -\kappa_V \sqrt{v \cdot V(\theta)}, \quad \forall \theta \in \Theta.$$

- ▶ If $\zeta(\theta) \leq \eta$, then $V(\theta) \leq \sqrt{\eta^{-2} + d} \cdot \sqrt{v}$. Therefore, letting $\alpha = \kappa_V / \sqrt{\eta^{-2} + d}$ we have $\nabla V(\theta)^T F(\theta) \leq -\alpha V(\theta)$.

Proposition 2

There exist constants $\eta > 0$ and $\beta > 0$ such that

$$\zeta(\theta) \geq \eta \implies \langle \nabla \zeta(\theta), F(\theta) \rangle \leq -\beta \cdot \zeta(\theta).$$

If $\zeta(\theta) < \eta$, then $\Psi_\eta(\theta) = \eta V(\theta)$ and by Proposition 1

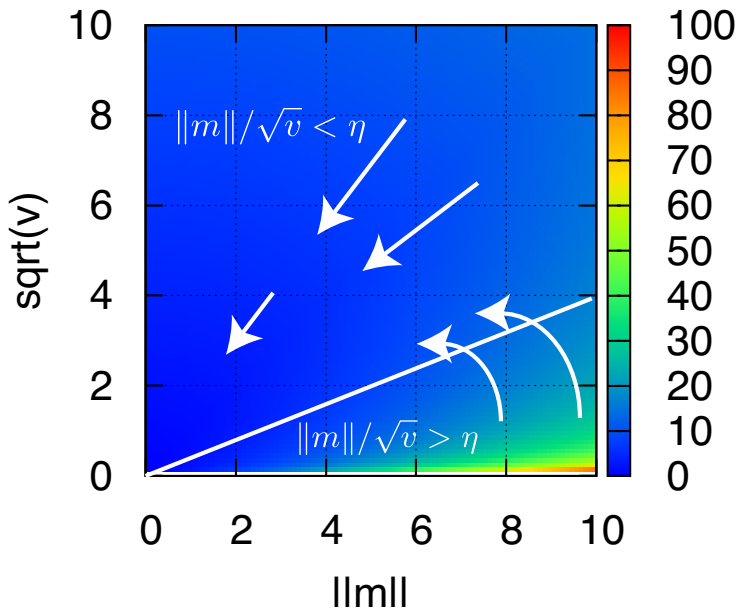
$$\begin{aligned}\nabla \Psi_\eta^2(\theta)^T F(\theta) &= \eta \nabla V(\theta)^T F(\theta) \\ &\leq -\alpha \eta V(\theta) = -\alpha \Psi_\eta^2(\theta) .\end{aligned}$$

If $\zeta(\theta) > \eta$, then $\Psi_\eta(\theta) = \zeta(\theta) V(\theta)$ by Proposition 2

$$\begin{aligned}\nabla \Psi_\eta^2(\theta)^T F(\theta) &= V(\theta) \cdot \underbrace{\nabla \zeta(\theta)^T F(\theta)}_{\leq -\beta \cdot \zeta(\theta)} + \zeta(\theta) \cdot \underbrace{\nabla V(\theta)^T F(\theta)}_{< 0} \\ &\leq -\beta \cdot \zeta(\theta) V(\theta) = -\beta \Psi_\eta^2(\theta) .\end{aligned}$$

Let $\kappa_e = \alpha \wedge \beta/2$, then

$$\begin{aligned}\frac{d}{dt} \Psi_\eta(\varphi(t, \theta)) &= \nabla \Psi_\eta(\varphi(t, \theta))^T F(\varphi(t, \theta)) \leq -\kappa_e \Psi_\eta(\varphi(t, \theta)) \\ \implies \Psi_\eta(\varphi(t, \theta)) &\leq \Psi_\eta(\theta) e^{-\kappa_e t} . \quad \square\end{aligned}$$



Proof of Proposition 1

Negative Correlation

Proposition 1

There exists a constant $\kappa_V > 0$ such that

$$\nabla V(\theta)^T F(\theta) \leq -\kappa_V \sqrt{v \cdot V(\theta)}, \quad \forall \theta \in \Theta.$$

$$\nabla V(\theta)^T F(\theta) \leq \frac{1}{\max_i A_{i,i}} \mathbb{E}_{x \sim P_\theta} \left[w(P_\theta[y : f(y) \leq f(x)]) (f(x) - \mathbb{E}_{x \sim P_\theta}[f(x)]) \right]$$

letting P_f be the distribution of $f(x)$

$$\leq \frac{1}{\max_i A_{i,i}} \mathbb{E}_f \left[\underbrace{w(P_f[y : y \leq f])}_{\text{nonincreasing w.r.t. } f} \overbrace{(f - \mathbb{E}_{f \sim P_f}[f])}^{\text{nondecreasing w.r.t. } f} \right]$$

Proof of Proposition 1

Sub propositions

Proposition 3 (negative correlation)

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$ be non-decreasing and non-increasing functions, respectively. Let μ be a measure on \mathbb{R} and suppose that $\int |f(x)| \mu(dx) < \infty$ and $\int |g(x)| \mu(dx) < \infty$. Let $M_f = \int f(x) \mu(dx)$ and $M_g = \int g(x) \mu(dx)$. Then,

$$\int f(x)g(x)\mu(dx) \leq M_f M_g - \frac{1}{4} \iint |f(x) - f(y)| \mu(dx)\mu(dy) \iint |g(x) - g(y)| \mu(dx)\mu(dy) .$$

Proposition 4 (fourth moment method)

Let Y and \tilde{Y} be i.i.d. random variable on \mathbb{R} such that $\mathbb{E}[Y^2] > 0$ and $\mathbb{E}[Y^4] < \infty$. Then, the following inequality holds

$$\mathbb{E}[|Y - \tilde{Y}|] \geq \frac{\mathbb{E}[(Y - \mathbb{E}[Y])^2]^{3/2}}{\mathbb{E}[(Y - \mathbb{E}[Y])^4]^{1/2}} .$$

Proof of Proposition 1

$$\nabla V(\theta)^T F(\theta) \leq \frac{1}{\max_i A_{i,i}} \mathbb{E}_f \left[\underbrace{w(P_f[y : y \leq f])}_{=: U \sim U(0,1)} \underbrace{(f - \mathbb{E}_{f \sim P_f}[f])}_{=: F} \right]$$

by the propositions

$$\leq -\frac{1}{4} \frac{1}{\max_i A_{i,i}} \frac{\mathbb{E}[(w(U) - \mathbb{E}[w(U)])^2]^{3/2}}{\mathbb{E}[(w(U) - \mathbb{E}[w(U)])^4]^{1/2}} \frac{\mathbb{E}[(F - \mathbb{E}[F])^2]^{3/2}}{\mathbb{E}[(F - \mathbb{E}[F])^4]^{1/2}}$$

since F is generalized χ^2 distributed, we can evaluate the expectations

$$= -\frac{1}{2\sqrt{30} \text{Cond}(A)} \frac{M_{w,2}^{3/2}}{M_{w,4}^{1/2}} \sqrt{v} \sqrt{2 \|m\|^2 + dv} \leq -\kappa_V \sqrt{vV(\theta)}$$

where $M_{w,1} := \int_0^1 w(u) du$, $M_{w,2} := \int_0^1 (w(u) - M_{w,1})^2 du$ and $M_{w,4} := \int_0^1 (w(u) - M_{w,1})^4 du$. □

Proof of Proposition 2

Proposition 2

There exist constants $\eta > 0$ and $\beta > 0$ such that

$$\zeta(\theta) \geq \eta \implies \nabla \zeta(\theta)^T F(\theta) \leq -\beta \cdot \zeta(\theta) .$$

Proof idea.

- ▶ $\frac{\nabla \zeta(\theta)^T F(\theta)}{\zeta(\theta)} = \frac{1}{\zeta(\theta)} \underbrace{\frac{m^T F^m(\theta)}{\|m\| v^{1/2}}}_{O(1)} - \frac{F^v(\theta)}{v}$
- ▶ $\frac{F^v(\theta)}{v} \rightarrow \frac{1}{d} \int w(P_1[y : y \leq z])(z^2 - 1)P_1(dz)$ as $\zeta(\theta) \rightarrow \infty$.
- ▶ By the assumption on w , $\int w(P_1[y : y \leq z])(z^2 - 1)P_1(dz) = \alpha > 0$.
- ▶ $\frac{\nabla \zeta(\theta)^T F(\theta)}{\zeta(\theta)} \rightarrow -\frac{\alpha}{d}$ as $\zeta(\theta) \rightarrow \infty$. □

Step 2: N -step cumulative error

1. Bound the progress of $\varphi(N\delta t, \theta_n)$ over $\theta_n = \varphi(0, \theta_n)$.

$$\Psi_\eta(\varphi(N\delta t, \theta_n)) \leq \Psi_\eta(\theta_n) e^{-\kappa_e N\delta t} .$$

2. Bound N -step cumulative error between θ_{n+N} and $\varphi(N\delta t, \theta_n)$.

3. Combine them and show $\mathbb{E}_n [\Psi_\eta(\theta_{N+n})] \leq \gamma \Psi_\eta(\theta_n)$.

N -step cumulative error

Lemma 2: N -step cumulative error

For any $T \geq 0$ and $n \in \mathbb{N}$, there is a constant $\kappa_T > 0$ independent of n such that

$$\mathbb{E}_n \left[\sup_{k \in \llbracket 0, T/\delta t \rrbracket} \rho(\theta_{n+k}, \varphi(k \cdot \delta t, \theta_n))^2 \right]^{1/2} \leq \kappa_T \sqrt{v_n} \sqrt{\delta t} .$$

Let $N = T/\delta t$. Then,

$$\mathbb{E}_n \left[\rho(\theta_{n+N}, \varphi(N\delta t, \theta_n))^2 \right]^{1/2} \leq \kappa_T \sqrt{v_n} \sqrt{\delta t} .$$

Proof Idea of Lemma 2

$$\begin{aligned}
 & \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \rho(\theta_{n+k}, \varphi(k\delta t, \theta_n)) \\
 & \leq \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left\| \sum_{i=0}^{k-1} \delta t (F^m(\theta_{n+i}) - F^m(\varphi(i\delta t, \theta_n))) \right\| \\
 & \quad + \frac{\sqrt{d}}{\sqrt{v_{n+1} \vee v^n(l\delta t)}} \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left| \sum_{i=0}^{k-1} \delta t (F^v(\theta_{n+i}) - F^v(\varphi(i\delta t, \theta_n))) \right| \\
 & \quad + \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left\| \sum_{i=0}^{k-1} \delta t M_{n+i}^m \right\| + \frac{\sqrt{d}}{\sqrt{v_{n+1} \vee v^n(l\delta t)}} \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left| \sum_{i=0}^{k-1} \delta t M_{n+i}^v \right| \\
 & \hspace{20em} \text{(martingale error)} \\
 & \quad + \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left\| \sum_{i=0}^{k-1} \int_{i\delta t}^{(i+1)\delta t} (F^m(\varphi(\tau, \theta_n)) - F^m(\varphi(i\delta t, \theta_n))) d\tau \right\| \\
 & \quad + \frac{\sqrt{d}}{\sqrt{v_{n+1} \vee v^n(l\delta t)}} \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left| \sum_{i=0}^{k-1} \int_{i\delta t}^{(i+1)\delta t} (F^v(\varphi(\tau, \theta_n)) - F^v(\varphi(i\delta t, \theta_n))) d\tau \right|. \\
 & \hspace{20em} \text{(discretization error)}
 \end{aligned}$$

Proof Idea of Lemma 2

Proposition: bound for first two terms

For some constant $K_1 > 0$,

$$\begin{aligned} & \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left\| \sum_{i=0}^{k-1} \delta t (F^m(\theta_{n+i}) - F^m(\varphi(i\delta t, \theta_n))) \right\| \\ & + \frac{\sqrt{d}}{\sqrt{v_{n+1} \vee v^n(l\delta t)}} \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left| \sum_{i=0}^{k-1} \delta t (F^v(\theta_{n+i}) - F^v(\varphi(i\delta t, \theta_n))) \right| \\ & \leq K_1 \delta t \sum_{i=0}^{\lfloor T/\delta t \rfloor - 1} \rho(\theta_{n+i}, \varphi(i\delta t, \theta_n)) . \end{aligned}$$

Proof Idea of Lemma 2

Proposition: bound for martingale error

For some constant $K_2 > 0$,

$$\begin{aligned} & \mathbb{E}_n \left[\left(\sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left\| \sum_{i=0}^{k-1} \delta t M_{n+i}^m \right\| \right)^2 \right]^{1/2} \\ & + \mathbb{E}_n \left[\left(\frac{\sqrt{d}}{\sqrt{v_{n+1}} \vee v^n(l\delta t)} \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left| \sum_{i=0}^{k-1} \delta t M_{n+i}^v \right| \right)^2 \right]^{1/2} \\ & \leq K_2 \sqrt{\delta t} \sqrt{v_n} \end{aligned}$$

Here we used Doob's L^p maximal inequality.

Proof Idea of Lemma 2

Proposition: bound for discretization error

For some constant $K_3 > 0$,

$$\begin{aligned} & \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left\| \sum_{i=0}^{k-1} \int_{i\delta t}^{(i+1)\delta t} (F^m(\varphi(\tau, \theta_n)) - F^m(\varphi(i\delta t, \theta_n))) d\tau \right\| \\ & + \frac{\sqrt{d}}{\sqrt{v_{n+1} \vee v^n(l\delta t)}} \sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \left\| \sum_{i=0}^{k-1} \int_{i\delta t}^{(i+1)\delta t} (F^v(\varphi(\tau, \theta_n)) - F^v(\varphi(i\delta t, \theta_n))) d\tau \right\| \\ & \leq K_3 \delta t \sqrt{v_n} \end{aligned}$$

Proof Idea of Lemma 2

By the propositions,

$$\begin{aligned}
 & \mathbb{E}_n \left[\sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \rho(\theta_{n+k}, \varphi(k\delta t, \theta_n))^2 \right]^{1/2} \\
 & \leq K_1 \delta t \sum_{i=0}^{\lfloor T/\delta t \rfloor - 1} \mathbb{E}_n \left[\rho(\theta_{n+i}, \varphi(i\delta t, \theta_n))^2 \right]^{1/2} + (K_2 + K_2 \sqrt{\delta t}) \sqrt{\delta t} \sqrt{v_n} \\
 & \leq K_1 \delta t \sum_{i=0}^{\lfloor T/\delta t \rfloor - 1} \mathbb{E}_n \left[\sup_{k \in \llbracket 0, i \rrbracket} \rho(\theta_{n+k}, \varphi(k\delta t, \theta_n))^2 \right]^{1/2} + (K_2 + K_3 \sqrt{\delta t}) \sqrt{\delta t} \sqrt{v_n}
 \end{aligned}$$

Apply discrete Gronwall inequality:

$$x_0 = 0, \quad x_{n+1} \leq C + \sum_{i=0}^n a_i x_i \implies x_{n+1} \leq C \exp \left(\sum_{i=0}^n a_i \right)$$

we have

$$\begin{aligned}
 \mathbb{E}_n \left[\sup_{k \in \llbracket 0, \lfloor T/\delta t \rrbracket \rrbracket} \rho(\theta_{n+k}, \varphi(k\delta t, \theta_n))^2 \right]^{1/2} & \leq (K_2 + K_2 \sqrt{\delta t}) \sqrt{\delta t} \sqrt{v_n} \exp(K_1 T) \\
 & \leq \kappa_T \sqrt{\delta t} \sqrt{v_n} \quad \square
 \end{aligned}$$

Step 3: Completion

1. Bound the progress of $\varphi(N\delta t, \theta_n)$ over $\theta_n = \varphi(0, \theta_n)$.

$$\Psi_\eta(\varphi(N\delta t, \theta_n)) \leq \Psi_\eta(\theta_n) e^{-\kappa_e N\delta t} .$$

2. Bound N -step cumulative error between θ_{n+N} and $\varphi(N\delta t, \theta_n)$.

$$\mathbb{E}_n \left[\rho(\theta_{n+N}, \varphi(N\delta t, \theta_n))^2 \right]^{1/2} \leq \kappa_T \sqrt{v_n} \sqrt{\delta t} .$$

3. Combine them and show $\mathbb{E}_n [\Psi_\eta(\theta_{N+n})] \leq \gamma \Psi_\eta(\theta_n)$.

Proof idea of Lemma 3

For simplicity, assume $\|m_n\| / \sqrt{v_n} \leq \eta$ for all n . Then,

$\Psi_\eta(\theta) = \eta^{1/2} \rho(\theta, \theta^*)$. We have

$$\begin{aligned}
 & \mathbb{E}_n[\Psi_\eta(\theta_{n+N})] \\
 &= \Psi_\eta(\varphi(N \cdot \delta t, \theta_n)) + \underbrace{\eta^{1/2} \cdot \mathbb{E}_n[(\rho(\theta_{n+N}, 0) - \rho(\varphi(N \cdot \delta t, \theta_n), 0))]}_{\text{triangle inequality}} \\
 &\leq \underbrace{\Psi_\eta(\varphi(N \cdot \delta t, \theta_n))}_{\text{Lemma 1}} + \underbrace{\eta^{1/2} \cdot \mathbb{E}_n[\rho(\theta_{n+N}, \varphi(N \cdot \delta t, \theta_n))]}_{\text{Lemma 2}} \\
 &\leq \Psi_\eta(\theta_n) e^{-\kappa_e N \delta t} + \kappa_T \sqrt{\delta t} \underbrace{\eta^{1/2} \sqrt{v_n}}_{\leq \Psi_\eta(\theta_n) / d^{1/2} \text{ by definition}} \\
 &= (e^{-\kappa_e N \delta t} + \kappa_T \sqrt{\delta t} / d^{1/2}) \Psi_\eta(\theta_n)
 \end{aligned}$$

taking sufficiently large $T = N \delta t$ and sufficiently small δt

$$\leq \gamma \Psi_\eta(\theta_n) .$$

$$\mathbb{E}_n[\Psi_\eta(\theta_{n+N})] \leq \gamma \Psi_\eta(\theta_n)$$

$$\implies \mathbb{E}[\Psi_\eta(\theta_{N \cdot k})] \leq \gamma^k \Psi_\eta(\theta_0)$$

$$\implies \frac{1}{k} \ln \mathbb{E}[\Psi_\eta(\theta_{N \cdot k})] \leq \ln \gamma + \frac{1}{k} \ln \Psi_\eta(\theta_0) \xrightarrow{k \rightarrow \infty} \ln \gamma < 0$$

$$\implies \limsup_{k \rightarrow \infty} \frac{1}{k} \ln \mathbb{E}[\Psi_\eta(\theta_{N \cdot k})] < 0$$

$$\implies \limsup_{k \rightarrow \infty} \frac{1}{k} \ln \mathbb{E} \left[\left(\|m_{N \cdot k} - \mathbf{x}^*\|^2 + d \cdot v_{N \cdot k} \right)^{1/2} \right] < 0$$



Global Linear Convergence in Expectation

on Monotonic Convex-Quadratic-Composite

Global Linear Convergence on Monotonic Convex-Quadratic-Composite

For sufficiently small $\delta t > 0$, there exists $\gamma < 0$ such that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \ln \mathbb{E} [\rho(\theta, \theta^*)] = \gamma .$$

Extension to Monotonic \mathcal{C}^2 -Composite

1. Bound the progress of $\varphi(N\delta t, \theta_n)$ over $\theta_n = \varphi(0, \theta_n)$.

Based on Taylor's approximation,

$$\Psi_\eta(\varphi(N\delta t, \theta_n)) \leq \Psi_\eta(\theta_n) e^{-\kappa_e N\delta t}, \quad \forall \theta_n \in U \subset \Theta .$$

2. Bound N -step cumulative error between θ_{n+N} and $\varphi(N\delta t, \theta_n)$.

(unchanged)

$$\mathbb{E}_n \left[\rho(\theta_{n+N}, \varphi(N\delta t, \theta_n))^2 \right]^{1/2} \leq \kappa_T \sqrt{v_n} \sqrt{\delta t} .$$

3. Combine them and we have

$$\Pr \left[\limsup_{n \rightarrow \infty} \frac{1}{n} \ln \Psi_\eta(\theta_n) < 0 \right] \geq 1 - \varepsilon .$$

With probability at most ε , θ_n will leave the region of attraction U .

Local Linear Convergence with probability $1 - \varepsilon$

on Monotonic \mathcal{C}^2 -Composite

Local Linear Convergence on Monotonic \mathcal{C}^2 -Composite

For $\theta_0 \in \mathcal{U}$ and any $0 < \varepsilon < 1$, there exists $\overline{\delta t} > 0$ such that for $0 < \forall \delta t \leq \overline{\delta t}$

$$\Pr \left[\limsup_{n \rightarrow \infty} \frac{1}{n} \ln \Psi_{\eta}(\theta_n, \theta^*) \leq \gamma \right] \geq 1 - \varepsilon ,$$

where the speed of convergence $\gamma < 0$ depends on δt .

Pros and Cons

Pros

- ▶ general weight scheme (cf $(1, \lambda)$ or $(\mu/\mu_I, \lambda)$ strategies)
- ▶ wide class of functions
- ▶ able to apply to general algorithms

Cons

- ▶ sufficiently small δt
 - ▶ the rate of convergence is loose
- ▶ rates of convergence of $\|m\|$ and \sqrt{v}
 - ▶ $\|m\|$ and \sqrt{v} converge at the same speed in practice and this is shown by [Auger 05] for $(1, \lambda)$ - σ SA-ES