

Convergence analysis of evolution strategies with covariance matrix adaptation

PhD students seminar

Armand Gissler

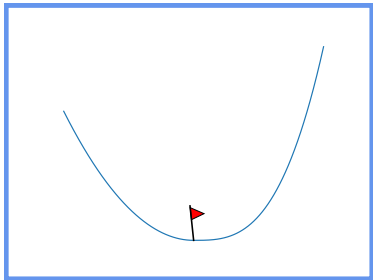
Wednesday 10th April, 2024

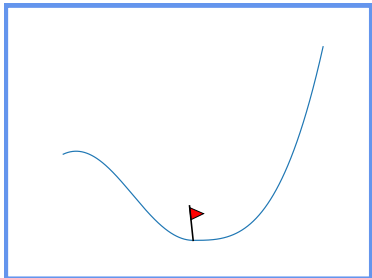
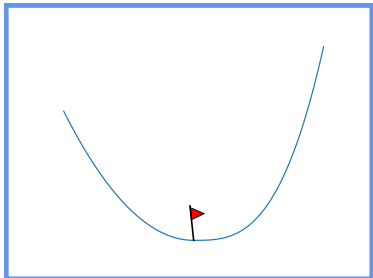
RandOpt team, Inria & École
polytechnique

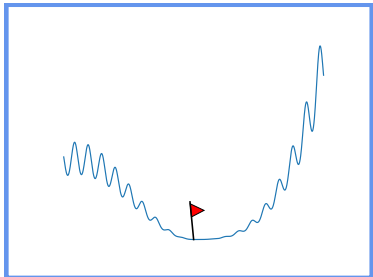
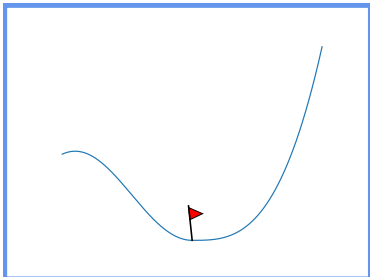
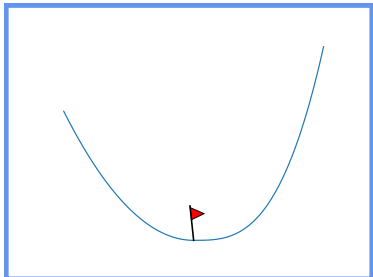
Advisors: Anne Auger & Nikolaus Hansen

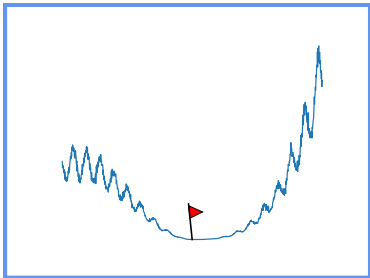
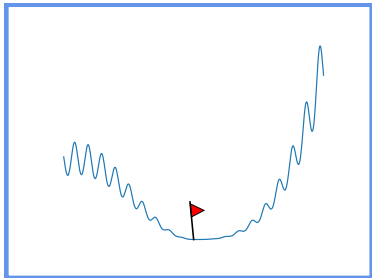
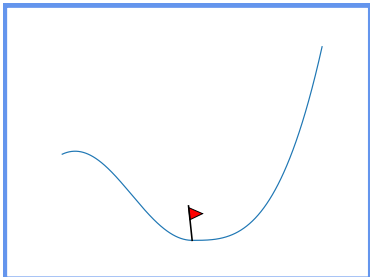
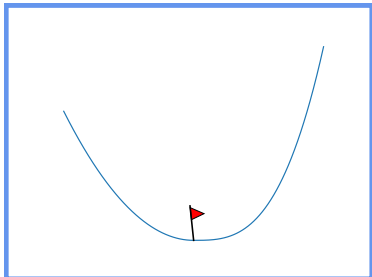
Inria

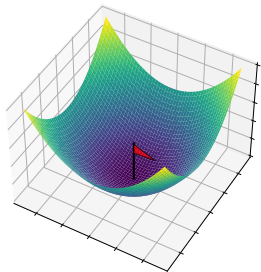


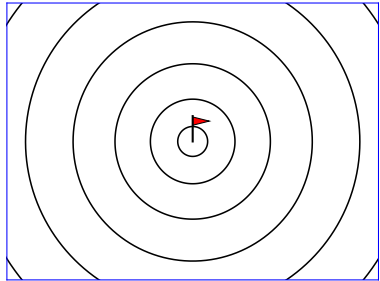
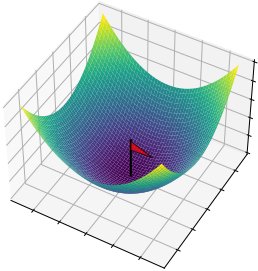


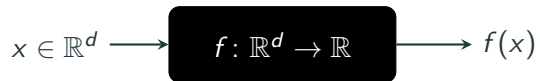


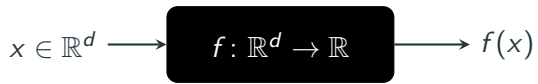






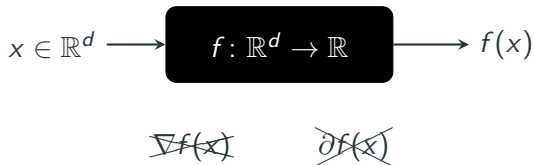




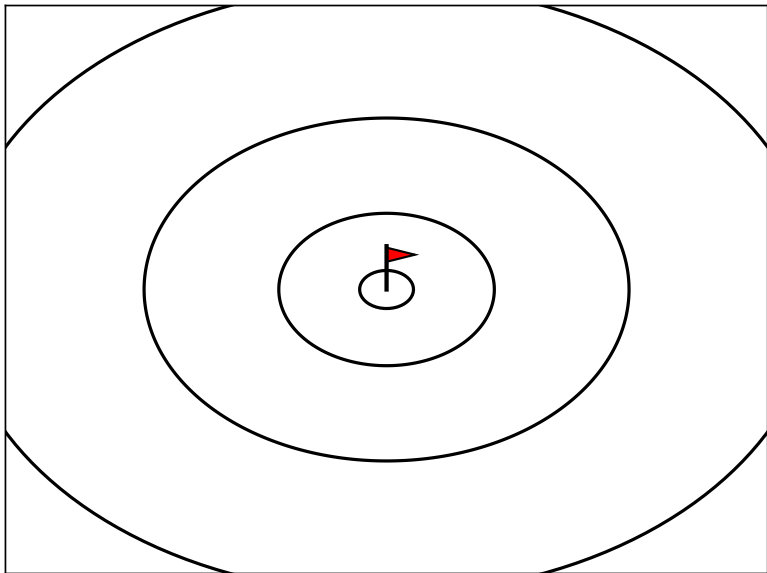


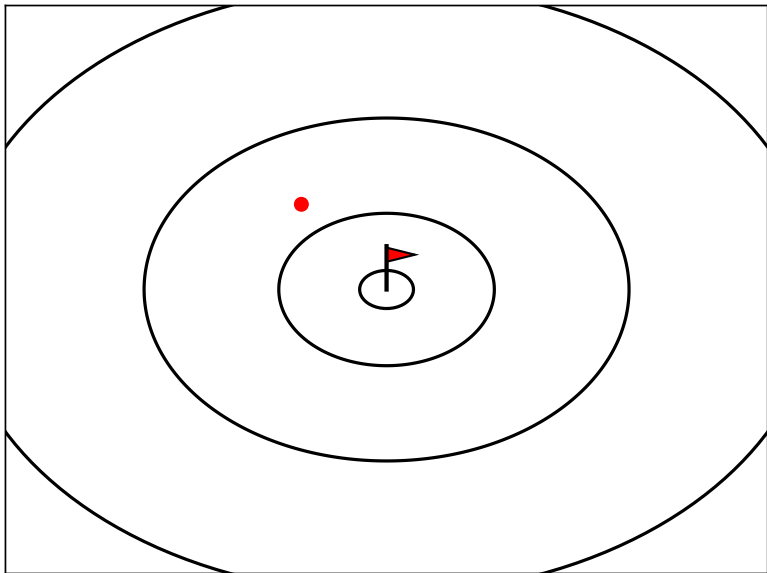
~~$\nabla f(x)$~~

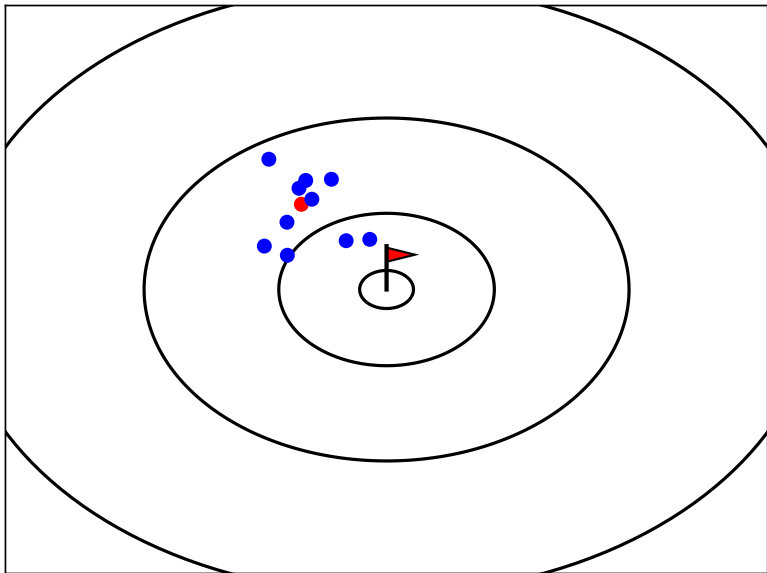
~~$\partial f(x)$~~

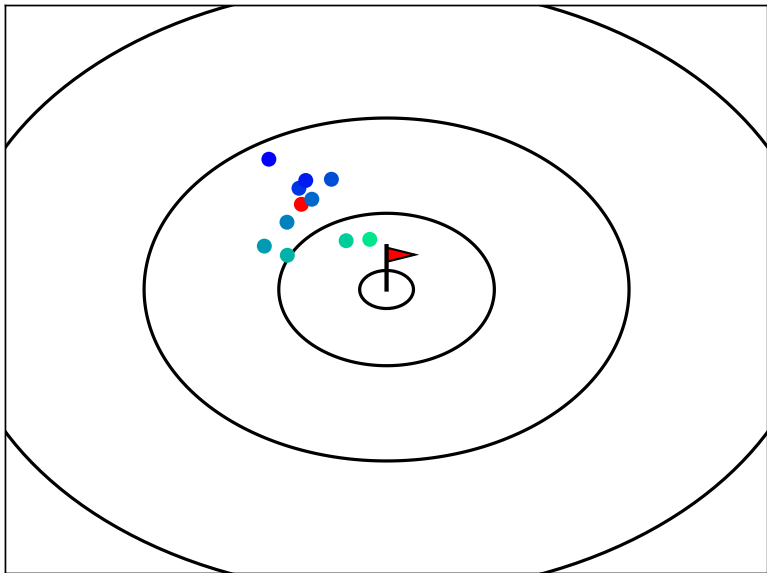


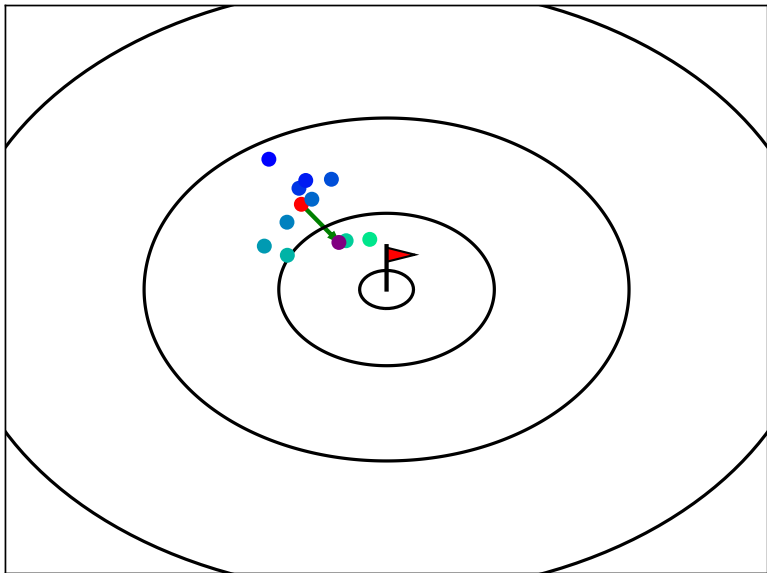
Find $x^* \in \underset{x \in \mathbb{R}^d}{\text{Arg min}} f(x)$ (P)





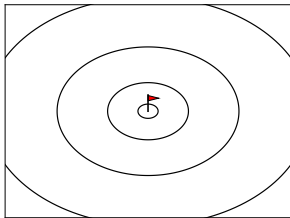






Algorithm 1 Our first ES

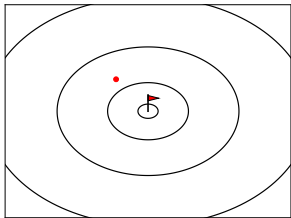
Goal: $\min_{x \in \mathbb{R}^d} f(x)$



Algorithm 1 Our first ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat: (Given $m_t \in \mathbb{R}^d$)

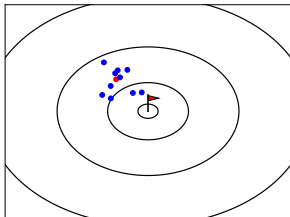


Algorithm 1 Our first ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat: (Given $m_t \in \mathbb{R}^d$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, I_d)$



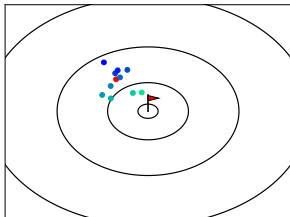
$\lambda =$ population size

Algorithm 1 Our first ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat: (Given $m_t \in \mathbb{R}^d$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, I_d)$
2. Rank population:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$



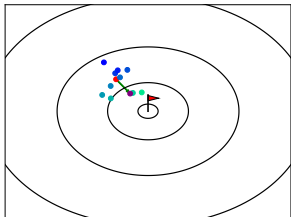
$\lambda =$ population size

Algorithm 1 Our first ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat: (Given $m_t \in \mathbb{R}^d$)

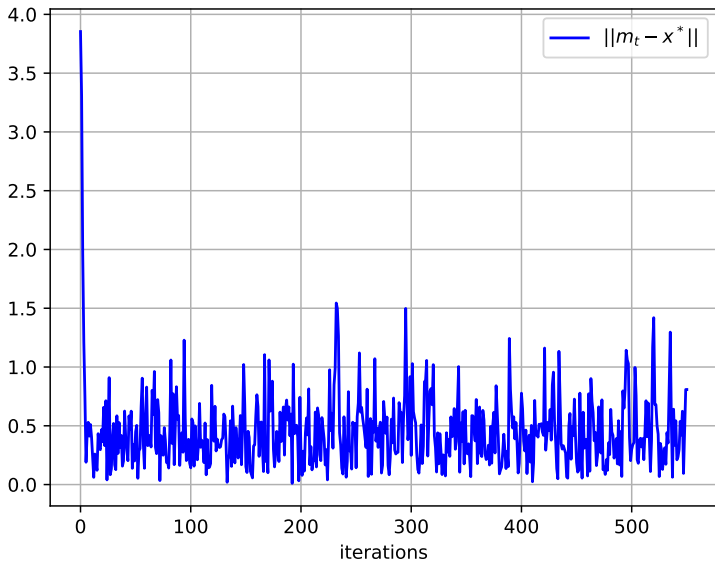
1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, I_d)$
2. Rank population:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update mean: $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$

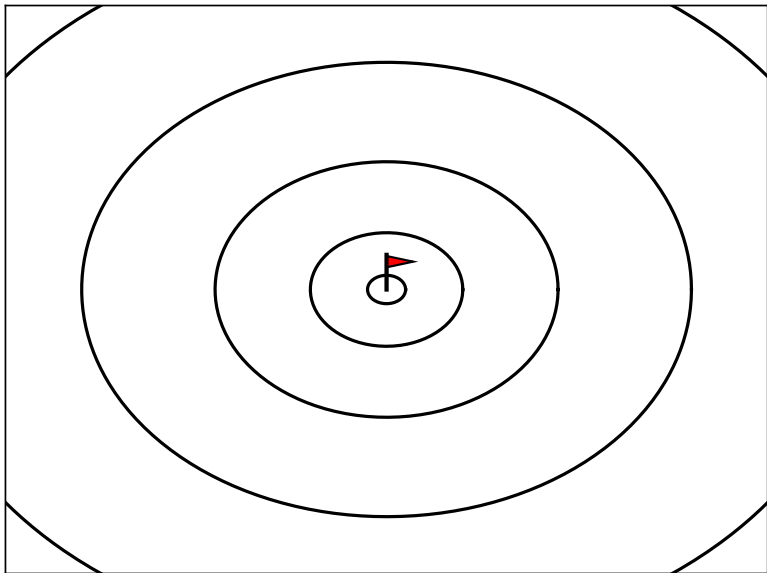


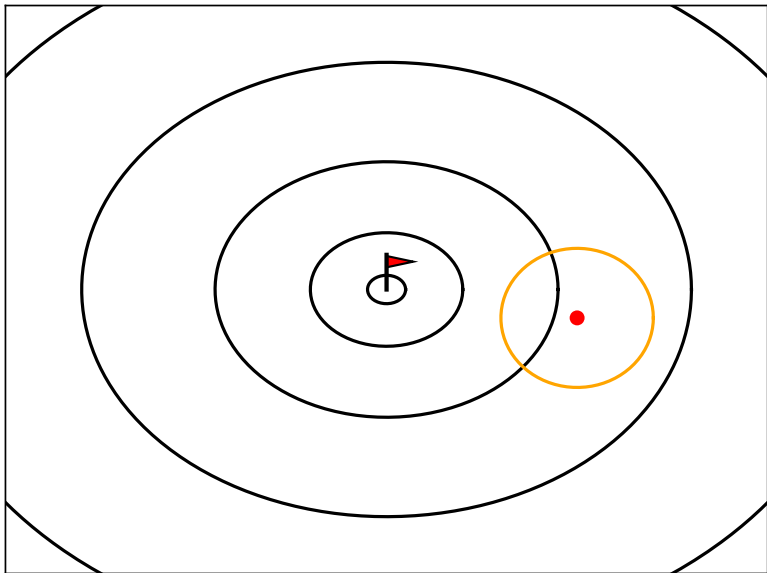
λ = population size

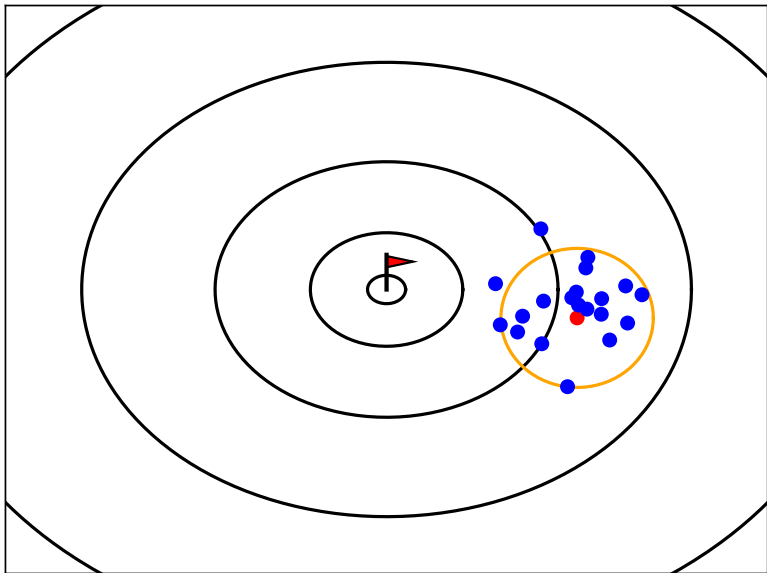
μ = parent number

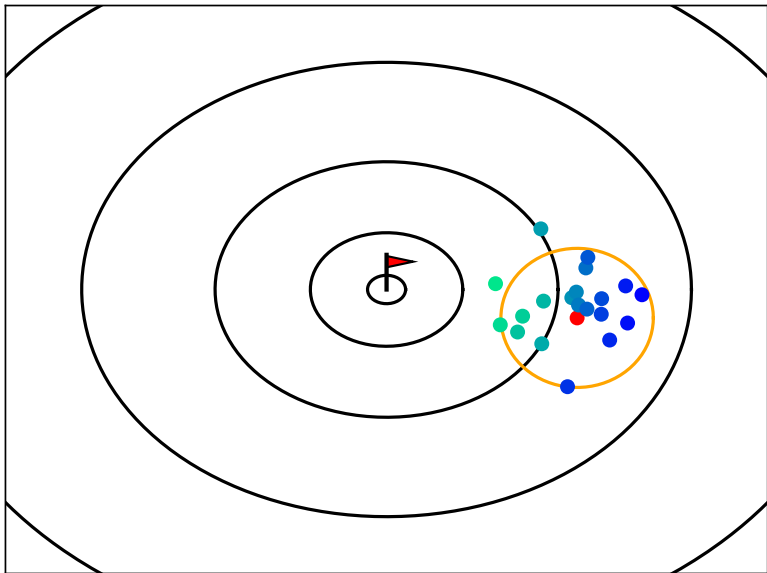
$$f: x \mapsto x^T A x$$

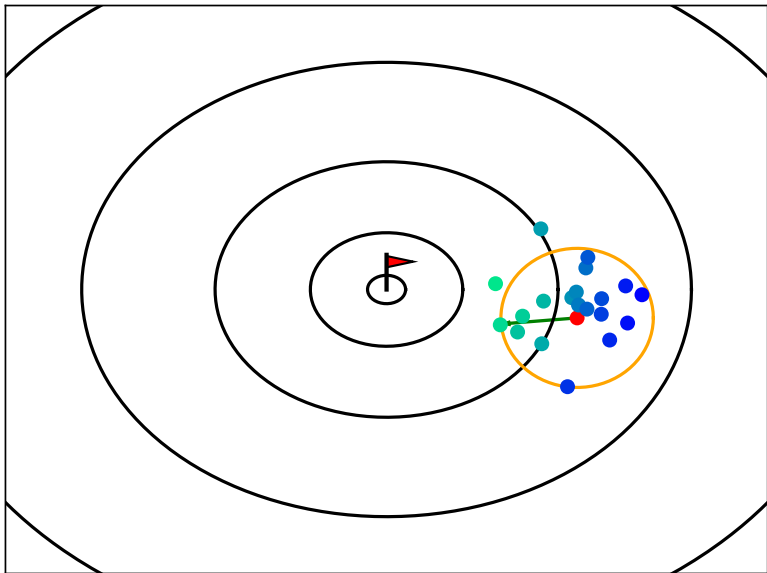


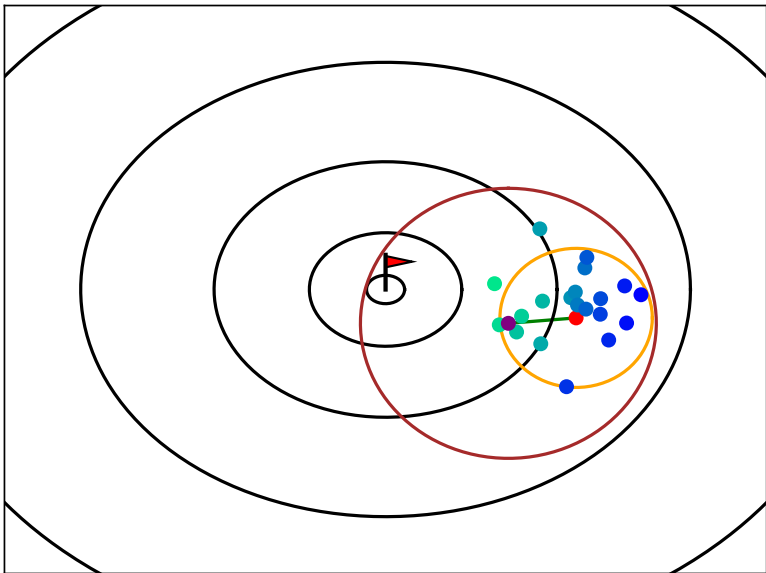


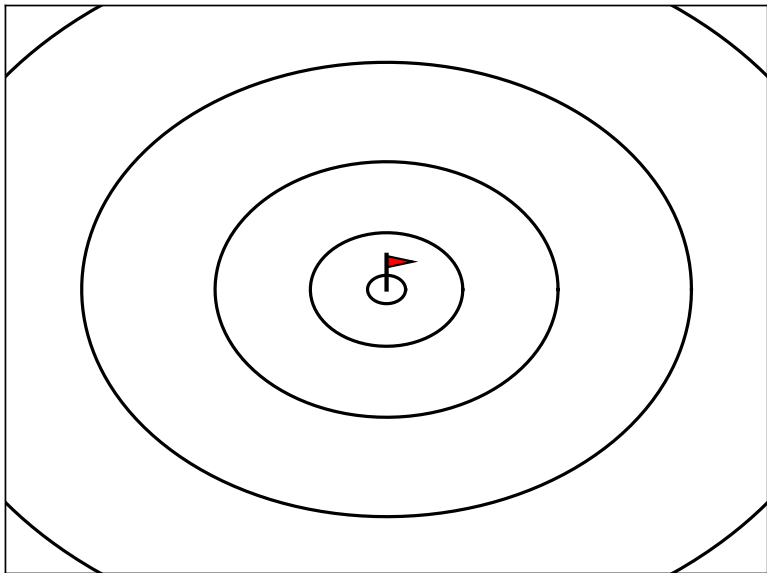


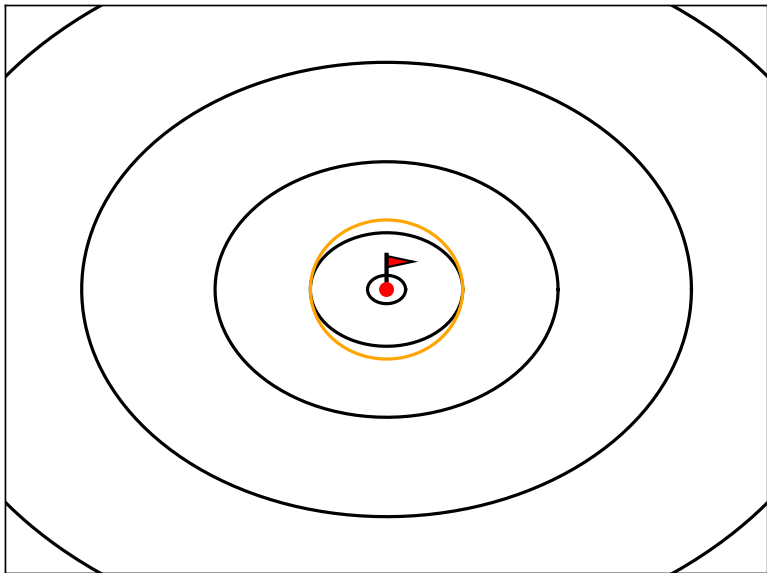


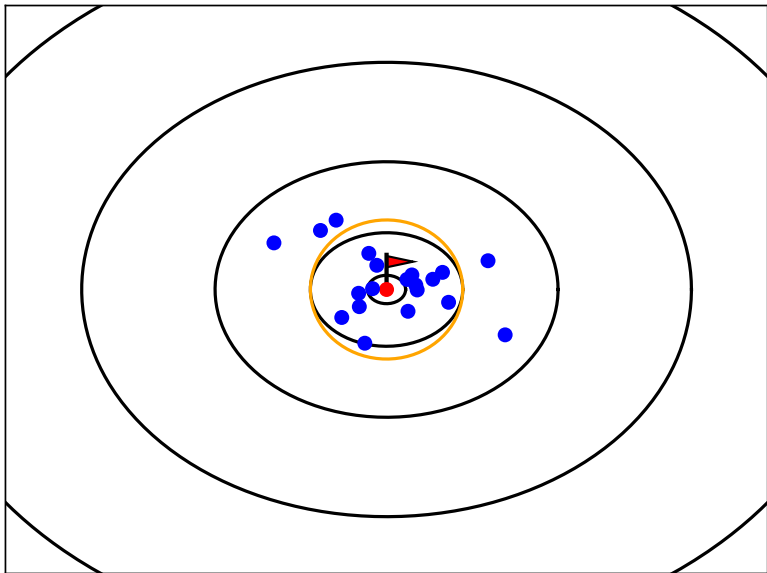


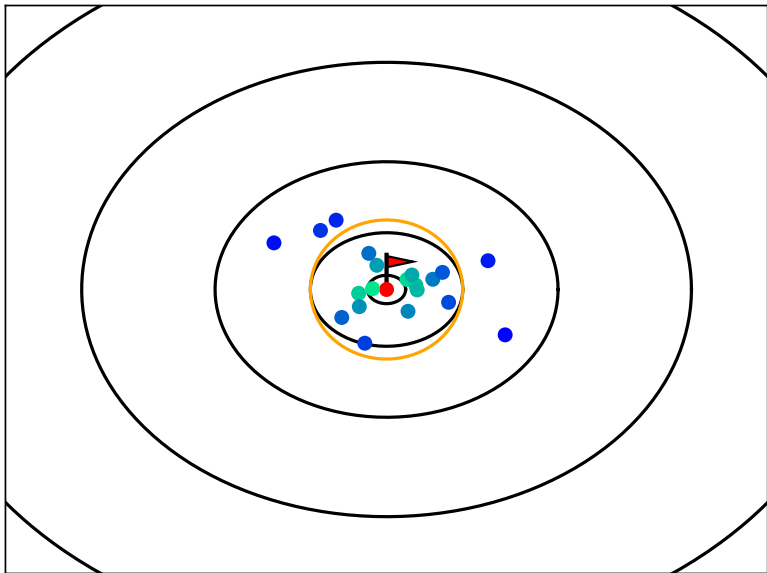


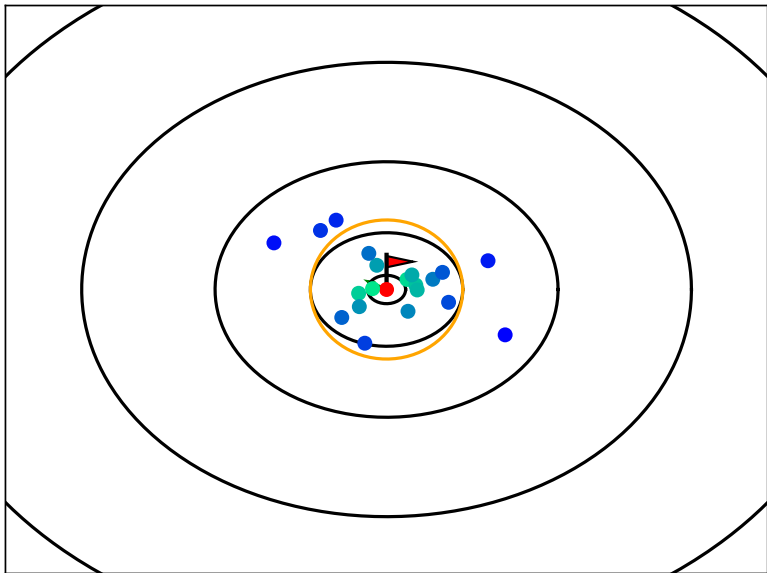


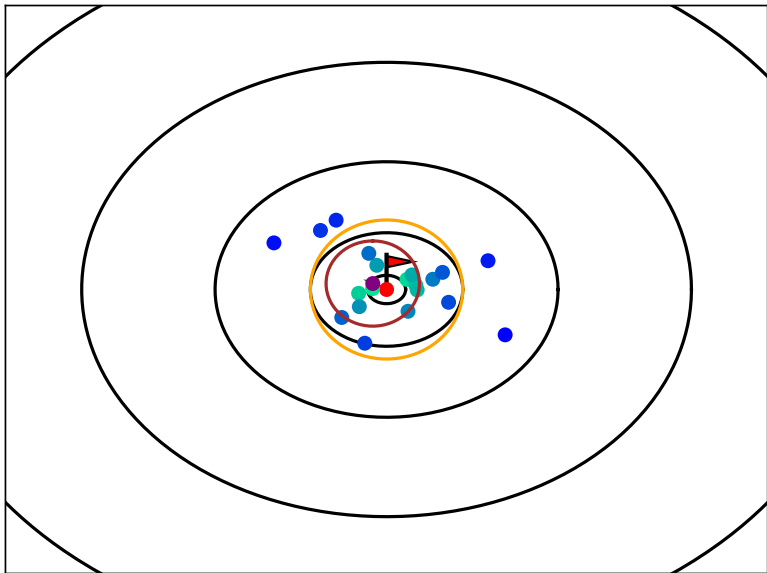






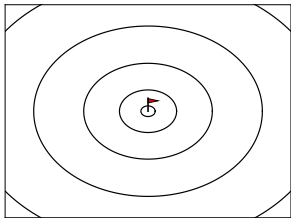






Algorithm 2 ES with step-size adaptation

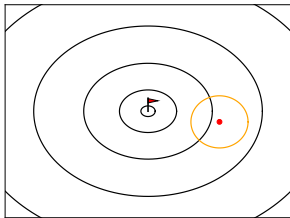
Goal: $\min_{x \in \mathbb{R}^d} f(x)$



Algorithm 2 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $m_t \in \mathbb{R}^d$ and $\sigma_t > 0$)

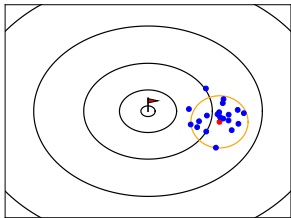


Algorithm 2 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $m_t \in \mathbb{R}^d$ and $\sigma_t > 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 I_d)$



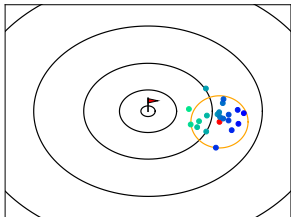
$\lambda =$ population size

Algorithm 2 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $m_t \in \mathbb{R}^d$ and $\sigma_t > 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 I_d)$
2. sort $f(x_{t+1}^i)$:



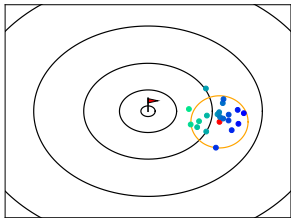
$\lambda =$ population size

Algorithm 2 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $m_t \in \mathbb{R}^d$ and $\sigma_t > 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 I_d)$
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$



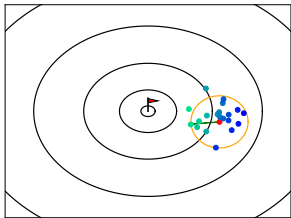
$\lambda =$ population size

Algorithm 2 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $m_t \in \mathbb{R}^d$ and $\sigma_t > 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 I_d)$
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$



λ = population size

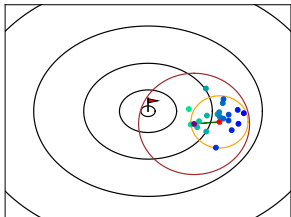
μ = parent number

Algorithm 2 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $m_t \in \mathbb{R}^d$ and $\sigma_t > 0$)

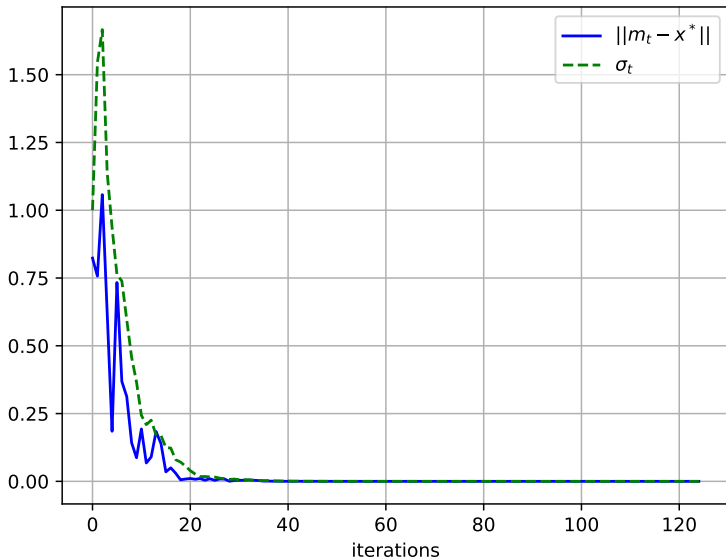
1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 I_d)$
 2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
 3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$
 4. $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$
-



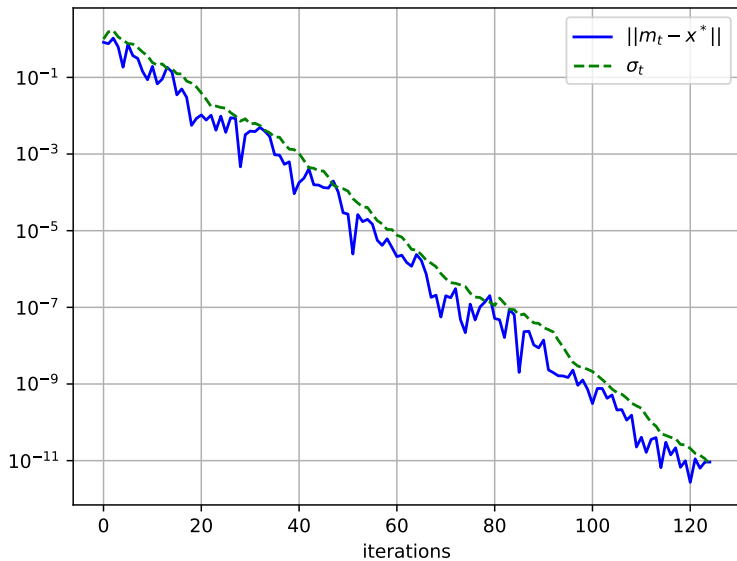
λ = population size

μ = parent number

$$f: x \mapsto x^T A x$$



$$f: x \mapsto x^T A x$$



Prove:

$$\frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \approx \frac{\sigma_{t+1}}{\sigma_t} \approx \rho \in (0, 1).$$

Prove:

$$\log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \approx \log \frac{\sigma_{t+1}}{\sigma_t} \approx -CR.$$

$P: X \times \mathcal{B}(X) \rightarrow \mathbb{R}$ is a **transition kernel** when

$\forall x \in X, P(x, \cdot)$ is a probability measure.

$P: X \times \mathcal{B}(X) \rightarrow \mathbb{R}$ is a **transition kernel** when

$\forall x \in X, P(x, \cdot)$ is a probability measure.

A **Markov chain** with transition kernel P is a random sequence $\{\theta_t\}_{t \in \mathbb{N}}$ such that:

$$\mathbb{P}[\theta_{t+1} \in A \mid \theta_t = x] = P(x, A).$$

$P: X \times \mathcal{B}(X) \rightarrow \mathbb{R}$ is a **transition kernel** when

$$\forall x \in X, \quad P(x, \cdot) \text{ is a probability measure.}$$

A **Markov chain** with transition kernel P is a random sequence $\{\theta_t\}_{t \in \mathbb{N}}$ such that:

$$\mathbb{P}[\theta_{t+1} \in A \mid \theta_t = x] = P(x, A).$$

- When $X = \{1, \dots, n\}$ is finite, P can be represented as a $n \times n$ matrix:

$$\mathbb{P}[\theta_{t+1} = j \mid \theta_t = i] = P_{ij}$$

$P: X \times \mathcal{B}(X) \rightarrow \mathbb{R}$ is a **transition kernel** when

$$\forall x \in X, \quad P(x, \cdot) \text{ is a probability measure.}$$

A **Markov chain** with transition kernel P is a random sequence $\{\theta_t\}_{t \in \mathbb{N}}$ such that:

$$\mathbb{P}[\theta_{t+1} \in A \mid \theta_t = x] = P(x, A).$$

- When $X = \{1, \dots, n\}$ is finite, P can be represented as a $n \times n$ matrix:

$$\mathbb{P}[\theta_{t+1} = j \mid \theta_t = i] = P_{ij}$$

- We can define a **k -steps transition kernel** P^k which satisfies

$$\mathbb{P}[\theta_{t+k} \in A \mid \theta_t = x] = P^k(x, A)$$

If X is finite:

If $X = \{1, \dots, n\}$:

$$\nu_0 = (p_1, \dots, p_n) \quad \text{with} \quad \sum_k p_k = 1$$

represents an initial state of the Markov chain $\{\theta_k\}_{k \in \mathbb{N}}$

If $X = \{1, \dots, n\}$:

$$\nu_0 = (p_1, \dots, p_n) \quad \text{with} \quad \sum_k p_k = 1$$

represents an initial state of the Markov chain $\{\theta_k\}_{k \in \mathbb{N}}$

After k steps:

$$\nu_k = \nu_0 \times P^k$$

If $X = \{1, \dots, n\}$:

$$\nu_0 = (p_1, \dots, p_n) \quad \text{with} \quad \sum_k p_k = 1$$

represents an initial state of the Markov chain $\{\theta_k\}_{k \in \mathbb{N}}$

After k steps:

$$\nu_k = \nu_0 \times P^k$$

If

$$\exists \pi, \forall \nu_0, \quad \lim_{k \rightarrow \infty} \nu_k = \pi$$

then $\{\theta_k\}_{k \in \mathbb{N}}$ is *ergodic*.

If X is **infinite**:

ν_0 **probability measure on X**

represents an initial state of the Markov chain $\{\theta_k\}_{k \in \mathbb{N}}$

After k steps:

$$\nu_k = \int \nu_0(dx) P^k(x, \cdot)$$

If

$$\exists \pi, \forall \nu_0, \quad \lim_{k \rightarrow \infty} \nu_k = \pi$$

then $\{\theta_k\}_{k \in \mathbb{N}}$ is *ergodic*.

For an ergodic Markov chain $\Theta = \{\theta_k\}_{k \in \mathbb{N}}$:

$$\theta_k \xrightarrow[k \rightarrow \infty]{} \pi$$

For an ergodic Markov chain $\Theta = \{\theta_k\}_{k \in \mathbb{N}}$:

$$\theta_k \xrightarrow[k \rightarrow \infty]{} \pi$$

where π is the invariant probability measure of Θ :

$$\theta_k \sim \pi \Rightarrow \theta_{k+1} \sim \pi$$

When $\Theta = \{\theta_k\}_{k \in \mathbb{N}}$ is ergodic:

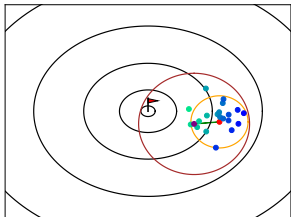
$$\lim_{T \rightarrow +\infty} \frac{1}{T} \sum_{k=0}^{T-1} g(\theta_k) = \int g(x) d\pi(x)$$

Algorithm 2 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $m_t \in \mathbb{R}^d$ and $\sigma_t > 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 I_d)$
 2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
 3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$
 4. $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$
-



λ = population size

μ = parent number

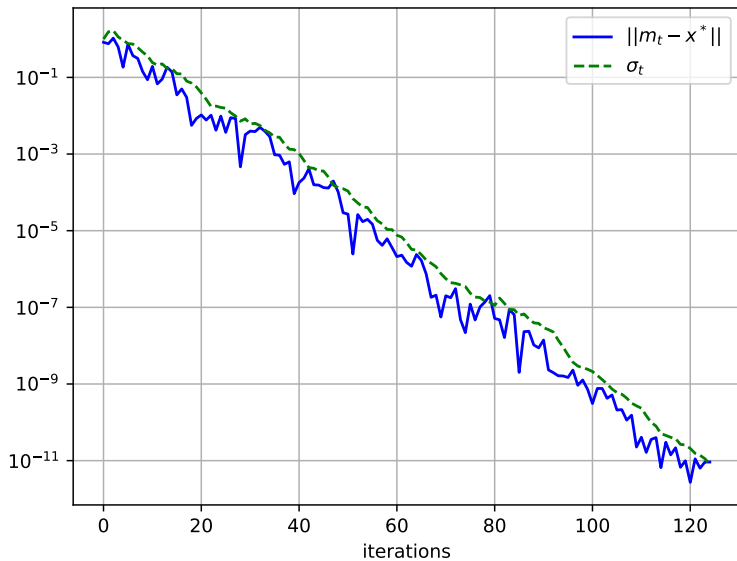
$\{(m_k, \sigma_k)\}_{k \in \mathbb{N}}$ is a Markov chain valued in $X = \mathbb{R}^d \times (0, +\infty)$

$$\lim_{k \rightarrow \infty} m_k = x^* \quad \text{and} \quad \lim_{k \rightarrow \infty} \sigma_k = 0$$

$$\lim_{k \rightarrow \infty} m_k = x^* \quad \text{and} \quad \lim_{k \rightarrow \infty} \sigma_k = 0$$

$\delta_{(x^*, 0)}$ is **not** a probability distribution on $X = \mathbb{R}^d \times (0, +\infty)$!

$$f: x \mapsto x^T A x$$



$$z_t = \frac{m_t - x^*}{\sigma_t}$$

$$z_t = \frac{m_t - \chi^*}{\sigma_t}$$

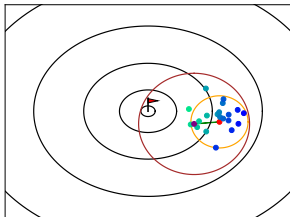
Question: $\{z_t\}_{t \in \mathbb{N}}$ is an ergodic Markov chain?

Algorithm 3 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $m_t \in \mathbb{R}^d$ and $\sigma_t > 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 I_d)$
 2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
 3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$
 4. $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$
-



λ = population size

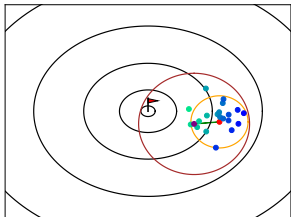
μ = parent number

Algorithm 3 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $z_t \in \mathbb{R}^d$) :

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 I_d)$
 2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
 3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$
 4. $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$
-



λ = population size

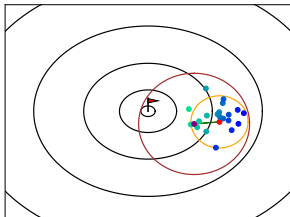
μ = parent number

Algorithm 3 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $z_t \in \mathbb{R}^d$) :

1. $z_{t+1}^1, \dots, z_{t+1}^\lambda \sim \mathcal{N}(z_t, l_d)$
 2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
 3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$
 4. $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$
-



λ = population size

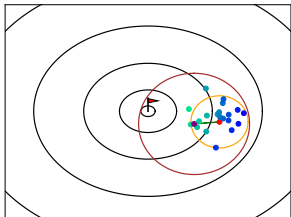
μ = parent number

Algorithm 3 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $z_t \in \mathbb{R}^d$) :

1. $z_{t+1}^1, \dots, z_{t+1}^\lambda \sim \mathcal{N}(z_t, l_d)$
 2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
 3. $z_{t+1} = \text{Average}(z_{t+1}^{1:\lambda}, \dots, z_{t+1}^{\mu:\lambda})$
 4. $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$
-



λ = population size

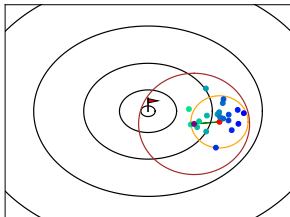
μ = parent number

Algorithm 3 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $z_t \in \mathbb{R}^d$) :

1. $z_{t+1}^1, \dots, z_{t+1}^\lambda \sim \mathcal{N}(z_t, I_d)$
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. $z_{t+1} = \frac{\text{Average}(z_{t+1}^{1:\lambda}, \dots, z_{t+1}^{\mu:\lambda})}{\text{increasing function}(\|z_{t+1} - z_t\|)}$



λ = population size

μ = parent number

Algorithm 3 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

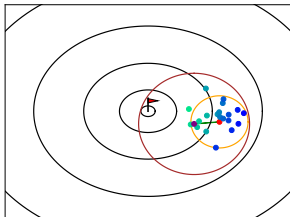
Repeat (Given $z_t \in \mathbb{R}^d$) :

1. $z_{t+1}^1, \dots, z_{t+1}^\lambda \sim \mathcal{N}(z_t, I_d)$

2. sort $f(x_{t+1}^i)$:

$$f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$$

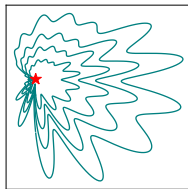
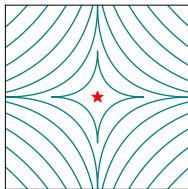
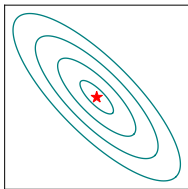
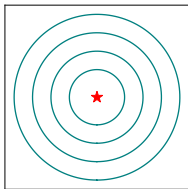
3. $z_{t+1} = \frac{\text{Average}(z_{t+1}^{1:\lambda}, \dots, z_{t+1}^{\mu:\lambda})}{\text{increasing function}(\|z_{t+1} - z_t\|)}$

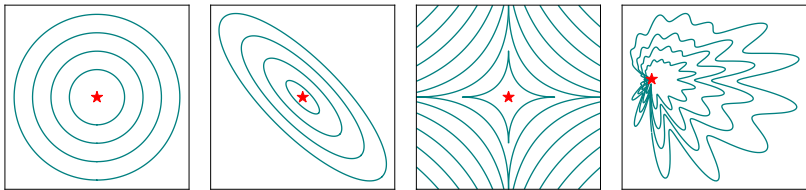


λ = population size

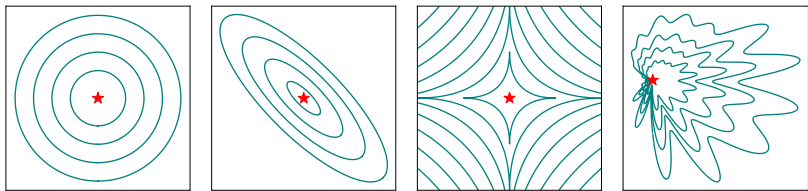
μ = parent number

$$f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda}) \stackrel{?}{\Leftrightarrow} g(z_{t+1}^{1:\lambda}) \leq \dots \leq g(z_{t+1}^{\lambda:\lambda})$$





$$f(m_t) \leq f(x_{t+1}) \Leftrightarrow f\left(\star + \frac{m_t - \star}{\sigma_t}\right) \leq f\left(\star + \frac{x_{t+1} - \star}{\sigma_t}\right)$$



$$f(m_t) \leq f(x_{t+1}) \Leftrightarrow f\left(\star + \frac{m_t - \star}{\sigma_t}\right) \leq f\left(\star + \frac{x_{t+1} - \star}{\sigma_t}\right)$$

Proposition

If $f \in \left\{ \left[\begin{array}{c} \text{concentric circles} \\ \text{elliptical contours} \\ \text{hyperbolic contours} \\ \text{jagged contours} \end{array} \right] \right\}$, then $\{z_t\}_{t \in \mathbb{N}}$ is a Markov chain.

$$z_t = \frac{m_t - \chi^*}{\sigma_t}$$

Question: $\{z_t\}_{t \in \mathbb{N}}$ is an ergodic **Markov chain**?

When X is finite:

Theorem

If $\{\theta_t\}_{t \in \mathbb{N}}$ is an irreducible, aperiodic Markov chain, then it is ergodic.

When X is finite:

Theorem

If $\{\theta_t\}_{t \in \mathbb{N}}$ is an irreducible, aperiodic Markov chain, then it is ergodic.

If $X = \{1, \dots, n\}$,

If $X = \{1, \dots, n\}$,

$\{\theta_t\}_{t \in \mathbb{N}}$ is irreducible if

$$\forall x, y \in X, \exists k > 0, P_{x,y}^k > 0.$$

If $X = \{1, \dots, n\}$,

$\{\theta_t\}_{t \in \mathbb{N}}$ is irreducible if

$$\forall x, y \in X, \underbrace{\exists k > 0, P_{x,y}^k > 0}_{x \rightsquigarrow y}.$$

If X infinite,

$\{\theta_t\}_{t \in \mathbb{N}}$ is irreducible if

$$\forall x \in X, A \subset X, \quad \text{volume}(A) > 0 \Rightarrow x \rightsquigarrow A$$

If X infinite,

$\{\theta_t\}_{t \in \mathbb{N}}$ is irreducible if

$$\forall x \in X, A \subset X, \quad \text{volume}(A) > 0 \Rightarrow x \rightsquigarrow A$$

for some volume on X .

When X is finite:

Theorem

If $\{\theta_t\}_{t \in \mathbb{N}}$ is an irreducible, aperiodic Markov chain, then it is ergodic.

If

$$\left\{ \begin{array}{l} \theta_1 \in A_1 \Rightarrow \mathbb{P}[\theta_2 \in A_2] = 1 \\ \theta_2 \in A_2 \Rightarrow \mathbb{P}[\theta_3 \in A_3] = 1 \\ \vdots \\ \theta_T \in A_T \Rightarrow \mathbb{P}[\theta_{T+1} \in A_1] = 1 \end{array} \right.$$

Then period = T .

When period = 1,

$\{\theta_t\}$ is **aperiodic**.

When X is finite:

Theorem

If $\{\theta_t\}_{t \in \mathbb{N}}$ is an irreducible, aperiodic Markov chain, then it is ergodic.

When X is infinite:

Theorem

If $\{\theta_t\}_{t \in \mathbb{N}}$ is an irreducible, aperiodic Markov chain, then it is ergodic if

$$\mathbb{E}[V(\theta_{t+1}) \mid \theta_t] \leq (1 - \varepsilon)V(\theta_t) \quad \text{if } \theta_t \notin \text{small set}$$

for some $V: X \rightarrow [1, +\infty]$.

Proposition: for $\{z_t\}_{t \in \mathbb{N}}$, compact sets are small

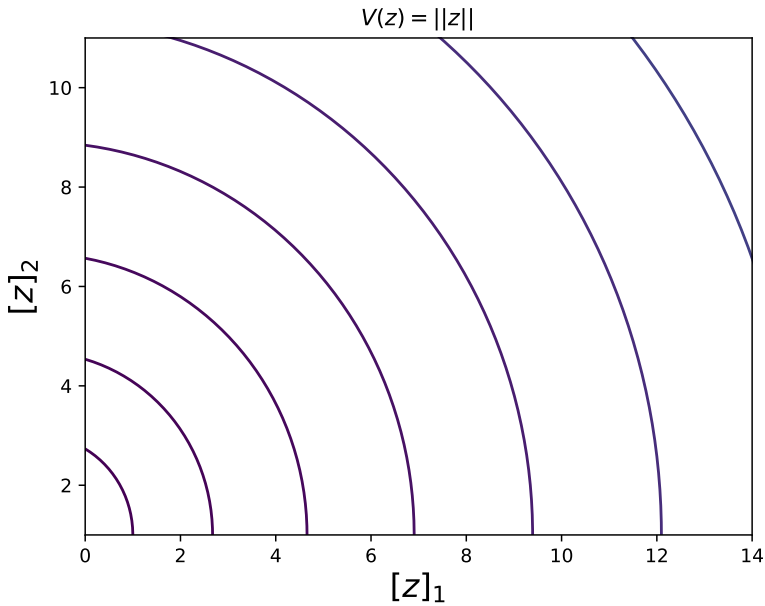
When X is infinite:

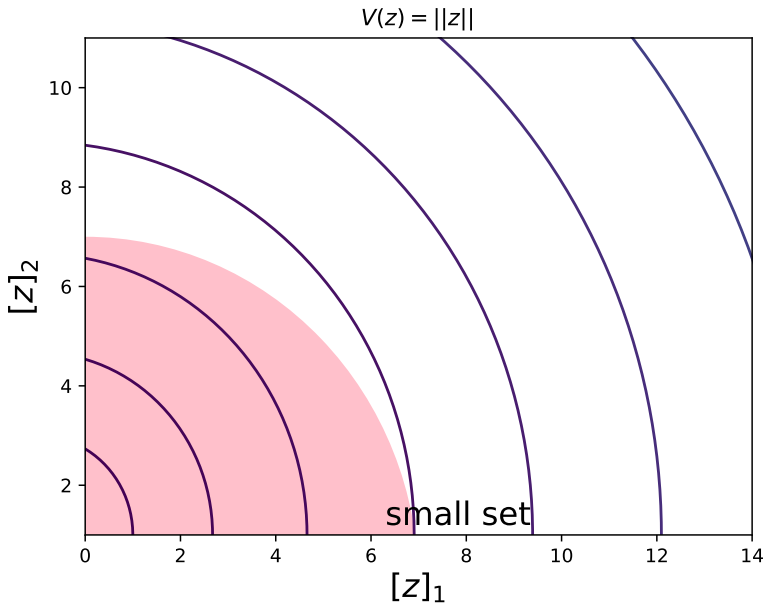
Theorem

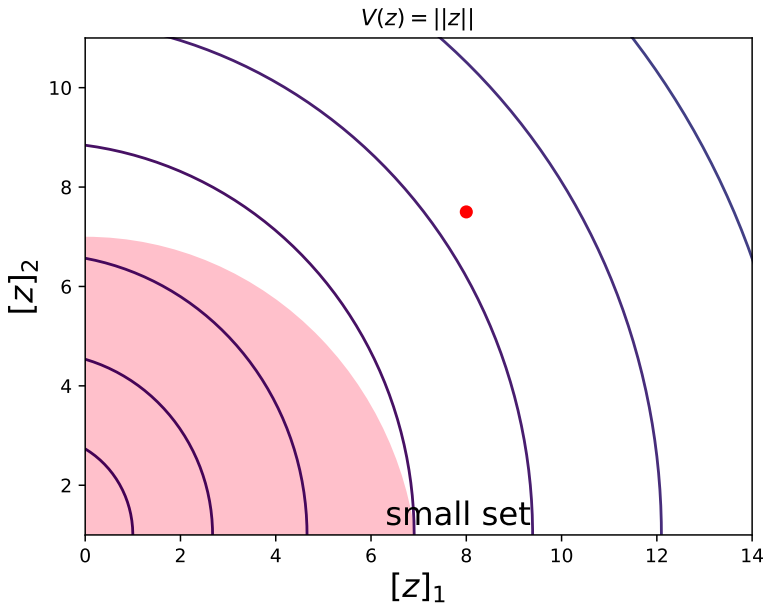
If $\{\theta_t\}_{t \in \mathbb{N}}$ is an irreducible, aperiodic Markov chain, then it is ergodic if

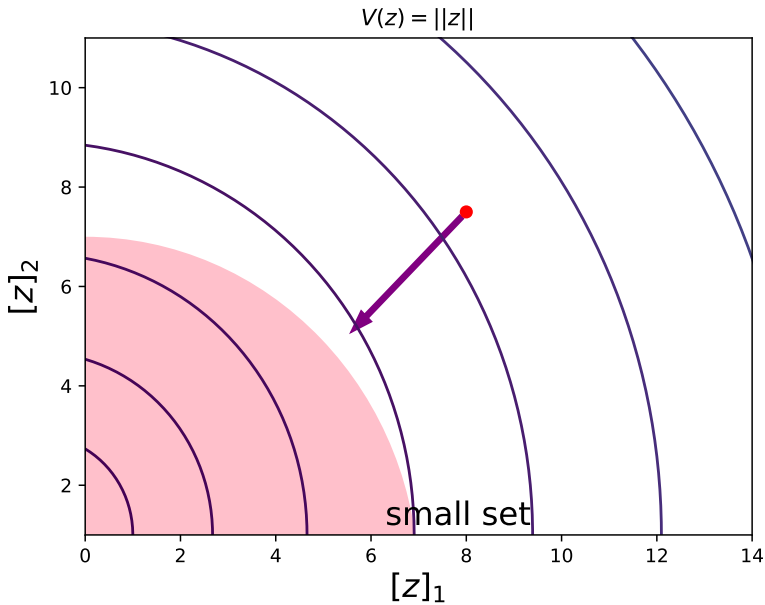
$$\mathbb{E}[V(\theta_{t+1}) \mid \theta_t] \leq (1 - \varepsilon)V(\theta_t) \quad \text{if } \theta_t \notin \text{compact set}$$

for some $V: X \rightarrow [1, +\infty]$.









When X is infinite:

Theorem

If $\{\theta_t\}_{t \in \mathbb{N}}$ is an irreducible, aperiodic Markov chain, then it is ergodic if

$$\mathbb{E}[V(\theta_{t+1}) \mid \theta_t] \leq (1 - \varepsilon)V(\theta_t) \quad \text{if } \theta_t \notin \text{compact set}$$

for some $V: X \rightarrow [1, +\infty]$.

Scheme of proof:

Scheme of proof:

1. irreducibility and aperiodicity of $\{z_t\}_{t \in \mathbb{N}}$

Scheme of proof:

1. irreducibility and aperiodicity of $\{z_t\}_{t \in \mathbb{N}}$
2. drift condition: $\exists K \subset \mathbb{R}^d$ compact and $V: \mathbb{R}^d \rightarrow [1, +\infty]$

$$\mathbb{E}[V(z_1)] \leq (1 - \varepsilon)V(z_0) \quad \forall z_0 \notin K$$

Scheme of proof:

1. irreducibility and aperiodicity of $\{z_t\}_{t \in \mathbb{N}}$
2. drift condition: $\exists K \subset \mathbb{R}^d$ compact and $V: \mathbb{R}^d \rightarrow [1, +\infty]$

$$\mathbb{E}[V(z_1)] \leq (1 - \varepsilon)V(z_0) \quad \forall z_0 \notin K$$

3. deduce convergence from the ergodicity

Scheme of proof:

1. **irreducibility and aperiodicity of $\{z_t\}_{t \in \mathbb{N}}$**
2. drift condition: $\exists K \subset \mathbb{R}^d$ compact and $V: \mathbb{R}^d \rightarrow [1, +\infty]$

$$\mathbb{E}[V(z_1)] \leq (1 - \varepsilon)V(z_0) \quad \forall z_0 \notin K$$

3. deduce convergence from the ergodicity

$$\theta_{k+1} = F(\theta_k, x_{k+1}) \quad (\mathbf{CM(F)})$$

where $x_{k+1} \sim p_{\theta_k}$

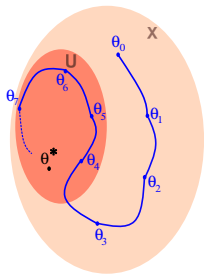
$$\begin{aligned}\theta_{k+1} &= F(\theta_k, x_{k+1}) && \text{(CM(F))} \\ &= F_{k+1}(\theta_0, x_1, \dots, x_{k+1})\end{aligned}$$

where $x_{k+1} \sim p_{\theta_k}$

θ^* is **attracting** when

θ^* is **attracting** when

$$\exists x_1, x_2, \dots, \lim_{k \rightarrow \infty} F_k(\theta_0, x_{1..k}) = \theta^*$$



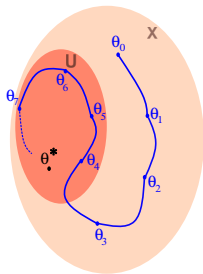
θ^* is **attracting** when

$$\exists x_1, x_2, \dots, \lim_{k \rightarrow \infty} F_k(\theta_0, x_{1..k}) = \theta^*$$

Theorem
If

- $\exists \theta^*$ *attracting*
- $\exists x_1^*, \dots, x_k^*$

such that $F_k(\theta^*, \cdot)$ is a **submersion** at $x_{1..k}^*$,
then, $\{\theta_t\}_{t \in \mathbb{N}}$ is **irreducible and aperiodic**.

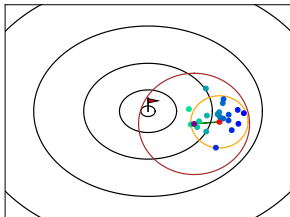


Algorithm 3 ES with step-size adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat (Given $z_t \in \mathbb{R}^d$) :

1. $z_{t+1}^1, \dots, z_{t+1}^\lambda \sim \mathcal{N}(z_t, I_d)$
2. sort $f(z_{t+1}^i)$:
 $f(z_{t+1}^{1:\lambda}) \leq \dots \leq f(z_{t+1}^{\lambda:\lambda})$
3. $z_{t+1} = \frac{\text{Average}(z_{t+1}^{1:\lambda}, \dots, z_{t+1}^{\mu:\lambda})}{\text{increasing function}(\|z_{t+1} - z_t\|)}$



λ = population size

μ = parent number

$$z_{k+1} = F(z_k, z_{k+1}^{1:\lambda}, \dots, z_{k+1}^{\lambda:\lambda})$$

$$z_{k+1} = F(z_k, z_{k+1}^{1:\lambda}, \dots, z_{k+1}^{\lambda:\lambda})$$

Proposition

0 is attracting

Proof.

Take $z_k^{i:\lambda} = 0$. Then

$$z_{k+1} = \frac{\text{Average}(0, \dots, 0)}{\text{normalization}} = 0$$

□

$$z_{k+1} = F(z_k, z_{k+1}^{1:\lambda}, \dots, z_{k+1}^{\lambda:\lambda})$$

Proposition

0 is attracting, and $F(0, \cdot)$ is submersive at 0.

Proof.

$$F(0, h^1, \dots, h^\lambda) = 0 + \text{Average}(h^1, \dots, h^\lambda) + o(h^1, \dots, h^\lambda)$$

□

Corollary

If $f \in \left\{ \begin{array}{c} \text{[concentric circles]} \\ \text{[elliptical contours]} \\ \text{[vertical wavy lines]} \\ \text{[chaotic fractal]} \end{array} \right\}$, $\{z_t\}_{t \in \mathbb{N}}$ is irreducible and aperiodic.

Scheme of proof:

1. irreducibility and aperiodicity of $\{z_t\}_{t \in \mathbb{N}}$
2. **drift condition:** $\exists K \subset \mathbb{R}^d$ **compact and** $V: \mathbb{R}^d \rightarrow [1, +\infty]$

$$\mathbb{E}[V(z_1)] \leq (1 - \varepsilon)V(z_0) \quad \forall z_0 \notin K$$

3. deduce convergence from the ergodicity

Proposition

If $f \in \left\{ \begin{array}{c} \text{[concentric circles]} \\ \text{[ellipses]} \\ \text{[chaotic attractor]} \end{array} \right\}$:

$$\mathbb{E}[\|z_{t+1}\| \mid z_t] \leq (1 - \varepsilon) \times \|z_t\| \quad \text{if } \|z_t\| \gg 1$$

Proposition

If $f \in \left\{ \begin{array}{c} \text{[concentric circles]} \\ \text{[ellipses]} \\ \text{[chaotic attractor]} \end{array} \right\}$:

$$\mathbb{E}[\|z_{t+1}\| \mid z_t] \leq (1 - \varepsilon) \times \|z_t\| \quad \text{if } \|z_t\| \gg 1$$

Theorem

If $f \in \left\{ \begin{array}{c} \text{[concentric circles]} \\ \text{[ellipses]} \\ \text{[chaotic attractor]} \end{array} \right\}$: $\{z_t\}_{t \in \mathbb{N}}$ is ergodic.

Scheme of proof:

1. irreducibility and aperiodicity of $\{z_t\}_{t \in \mathbb{N}}$
2. drift condition: $\exists K \subset \mathbb{R}^d$ compact and $V: \mathbb{R}^d \rightarrow [1, +\infty]$

$$\mathbb{E}[V(z_1)] \leq (1 - \varepsilon)V(z_0) \quad \forall z_0 \notin K$$

3. **deduce convergence from the ergodicity**

Theorem

If $f \in \left\{ \begin{array}{c} \text{[concentric circles]} \\ \text{[ellipses]} \\ \text{[chaotic fractal]} \end{array} \right\}$, ES converges linearly (or geometrically):

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = \lim_{t \rightarrow \infty} \mathbb{E} \left[\log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \right] = -\text{CR}.$$

Proof.

$$\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|}$$

□

Theorem

If $f \in \left\{ \begin{array}{c} \text{[concentric circles]} \\ \text{[ellipses]} \\ \text{[chaotic fractal]} \end{array} \right\}$, ES converges linearly (or geometrically):

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = \lim_{t \rightarrow \infty} \mathbb{E} \left[\log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \right] = -\text{CR}.$$

Proof.

$$\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} = \frac{1}{T} \sum_{t=0}^{T-1} \log \|m_{t+1} - x^*\| - \log \|m_t - x^*\|$$

□

Theorem

If $f \in \left\{ \begin{array}{c} \text{[Target Function]} \\ \text{[Convex Function]} \\ \text{[Non-convex Function]} \end{array} \right\}$, ES converges linearly (or geometrically):

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = \lim_{t \rightarrow \infty} \mathbb{E} \left[\log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \right] = -\text{CR}.$$

Proof.

$$\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} = \frac{1}{T} \sum_{t=0}^{T-1} \log \sigma_{t+1} \|z_{t+1}\| - \log \sigma_t \|z_t\|$$

□

Theorem

If $f \in \left\{ \begin{array}{c} \text{[concentric circles]} \\ \text{[ellipses]} \\ \text{[chaotic fractal]} \end{array} \right\}$, ES converges linearly (or geometrically):

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = \lim_{t \rightarrow \infty} \mathbb{E} \left[\log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \right] = -\text{CR}.$$

Proof.

$$\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} = \frac{1}{T} \sum_{t=0}^{T-1} \log \|z_{t+1}\| - \log \|z_t\| + \log \frac{\sigma_{t+1}}{\sigma_t}$$

□

Theorem

If $f \in \left\{ \begin{array}{c} \text{[concentric circles]} \\ \text{[ellipses]} \\ \text{[fractal]} \end{array} \right\}$, ES converges linearly (or geometrically):

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = \lim_{t \rightarrow \infty} \mathbb{E} \left[\log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \right] = -\text{CR}.$$

Proof.

$$\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} = \frac{1}{T} \sum_{t=0}^{T-1} \log \|z_{t+1}\| - \log \|z_t\| + \log \uparrow (\|z_{t+1} - z_t\|)$$

□

Theorem

If $f \in \left\{ \begin{array}{c} \text{[Target Function]} \\ \text{[Elliptical Level Sets]} \\ \text{[Convex Level Sets]} \end{array} \right\}$, ES converges linearly (or geometrically):

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = \lim_{t \rightarrow \infty} \mathbb{E} \left[\log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \right] = -\text{CR}.$$

Proof.

$$\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} = \mathbb{E}_{z_t \sim \pi} [\log \|z_{t+1}\| - \log \|z_t\| + \uparrow (\|z_{t+1} - z_t\|)]$$

□

Theorem

If $f \in \left\{ \begin{array}{c} \text{[concentric circles]} \\ \text{[ellipses]} \\ \text{[chaotic fractal]} \end{array} \right\}$, ES converges linearly (or geometrically):

$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = \lim_{t \rightarrow \infty} \mathbb{E} \left[\log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \right] = -\text{CR}.$$

Proof.

$$\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} = \mathbb{E}_{z_t \sim \pi} [\uparrow (\|z_{t+1} - z_t\|)]$$

□

Theorem

If $f \in \left\{ \begin{array}{c} \text{[target]}, \text{[contour]}, \text{[stochastic]} \end{array} \right\}$, ES converges linearly (or geometrically):

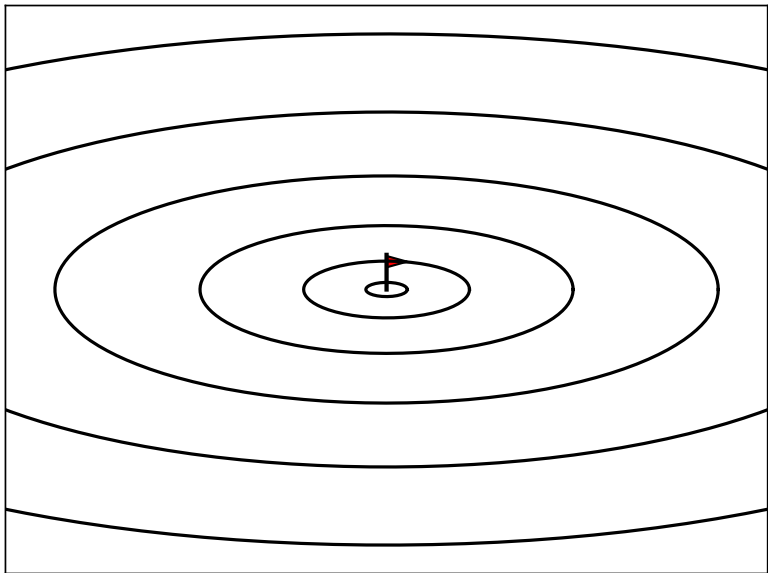
$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = \lim_{t \rightarrow \infty} \mathbb{E} \left[\log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \right] = -\text{CR}.$$

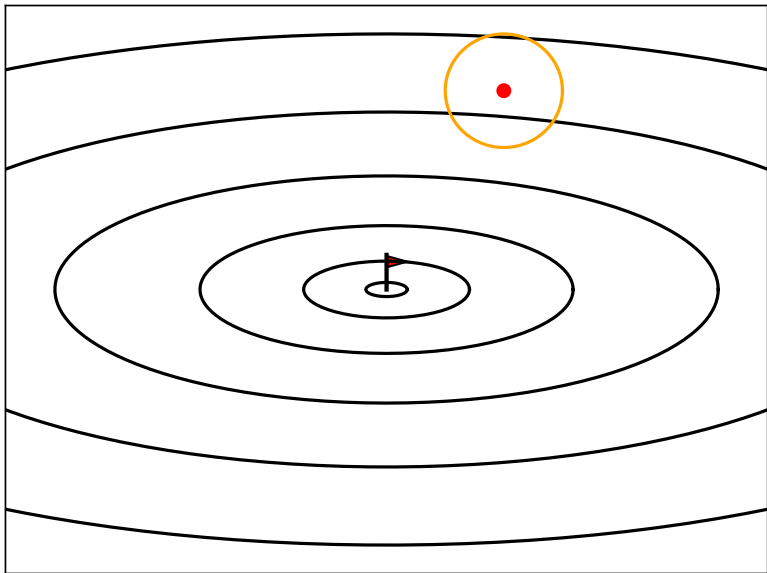
Proof.

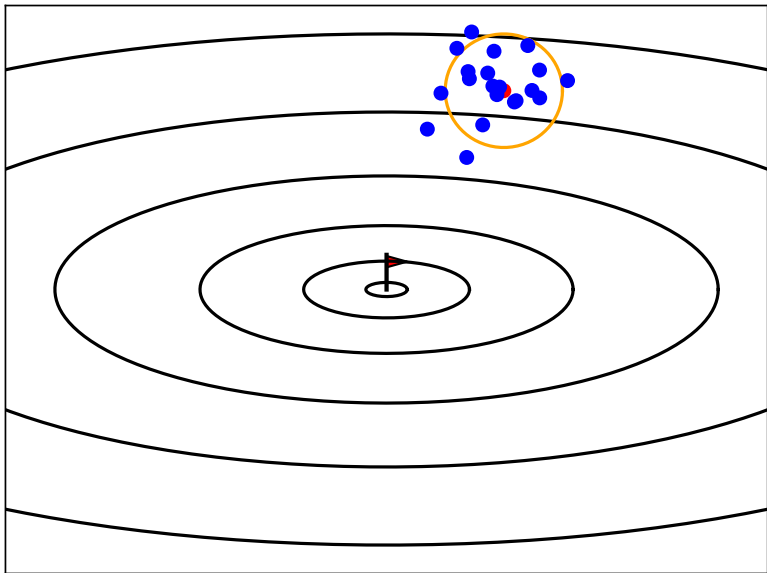
$$\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} = \mathbb{E}_{z_t \sim \pi} [\uparrow (\|z_{t+1} - z_t\|)]$$

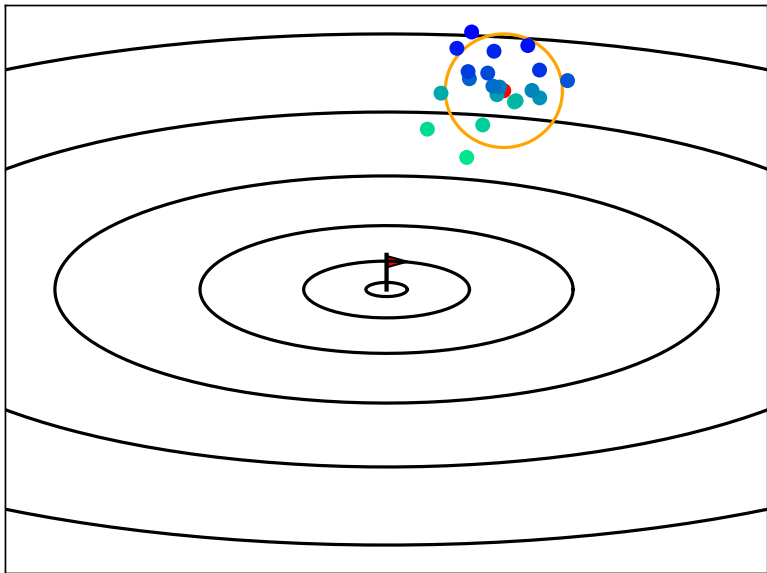
$$\text{CR} = -\mathbb{E}_{z_t \sim \pi} \uparrow (\|z_{t+1} - z_t\|)$$

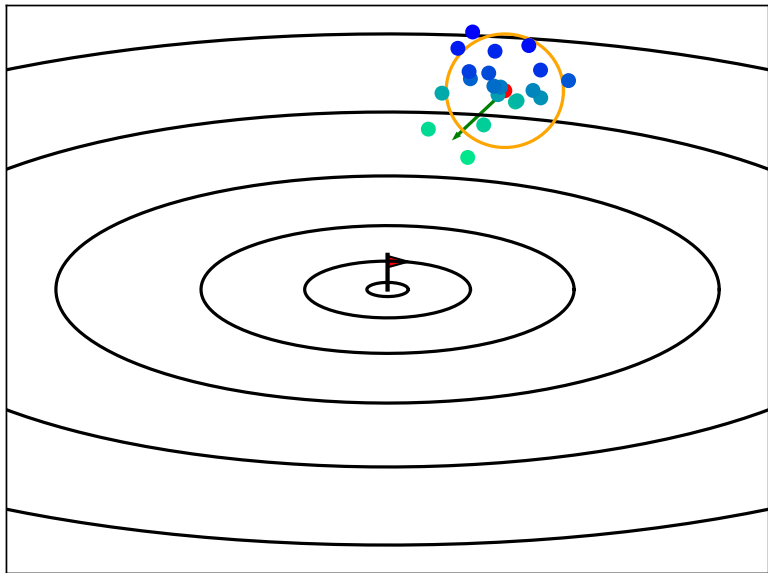
□

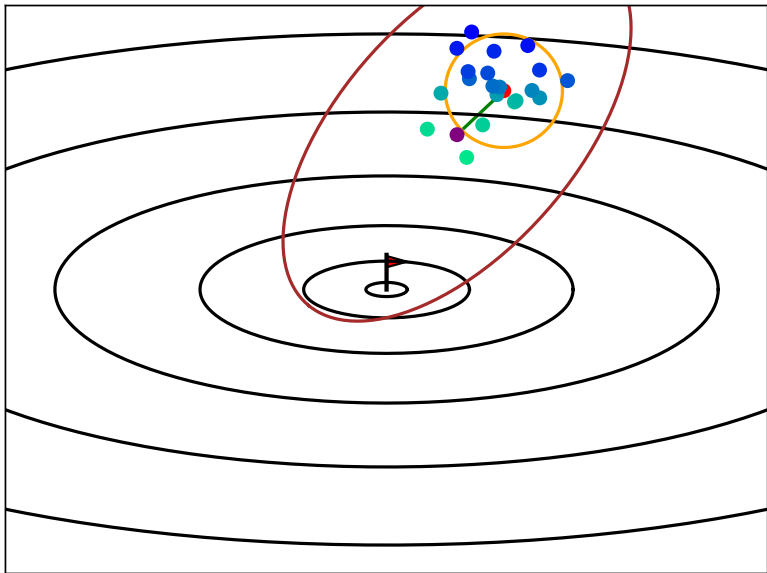












$$x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$$

where C_t is the covariance matrix at iteration t

$$x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$$

where C_t is the covariance matrix at iteration t

$x_{t+1}^{1:\lambda}$ is the best point among $\{x_{t+1}^1, \dots, x_{t+1}^\lambda\}$

$$x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$$

where C_t is the covariance matrix at iteration t

$x_{t+1}^{1:\lambda}$ is the best point among $\{x_{t+1}^1, \dots, x_{t+1}^\lambda\}$

Idea: sample more in the direction $x_{t+1}^{1:\lambda} - m_t$ at iteration $t + 1$

If v is a vector,

\overleftrightarrow{v} is a matrix with range collinear to v

If v is a vector,

\overleftrightarrow{v} is a matrix with range collinear to v :

$$\overleftrightarrow{v} = v \otimes v = v \times v^T$$

If v is a vector,

\overleftrightarrow{v} is a matrix with range collinear to v :

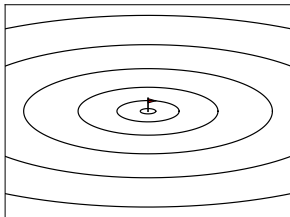
$$\overleftrightarrow{v} = v \otimes v = v \times v^T$$

$$C_{t+1} = \text{Positive combination} \left(C_t, \overleftrightarrow{x_{t+1}^{1:\lambda} - m_t} \right)$$

favors more the sampling in the direction $x_{t+1}^{1:\lambda} - m_t$ than C_t

Algorithm 4 ES with covariance matrix adaptation

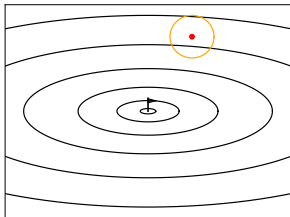
Goal: $\min_{x \in \mathbb{R}^d} f(x)$



Algorithm 4 ES with covariance matrix adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

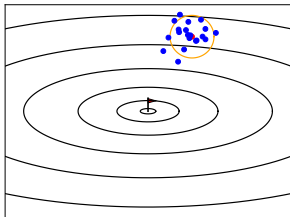


Algorithm 4 ES with covariance matrix adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$



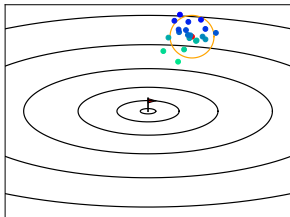
$\lambda =$ population size

Algorithm 4 ES with covariance matrix adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$
2. sort $f(x_{t+1}^i)$:



$\lambda =$ population size

Algorithm 4 ES with covariance matrix adaptation

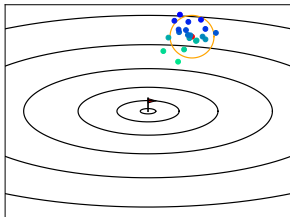
Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$

2. sort $f(x_{t+1}^i)$:

$$f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$$



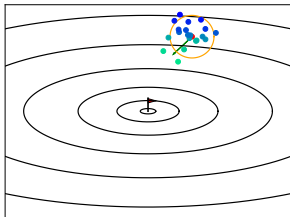
$\lambda =$ population size

Algorithm 4 ES with covariance matrix adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$



λ = population size

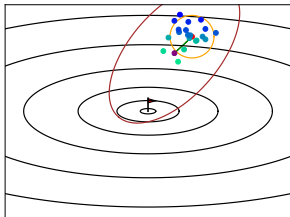
μ = parent number

Algorithm 4 ES with covariance matrix adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$
4. $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$



λ = population size

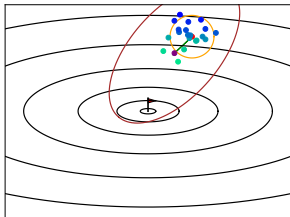
μ = parent number

Algorithm 4 ES with covariance matrix adaptation

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

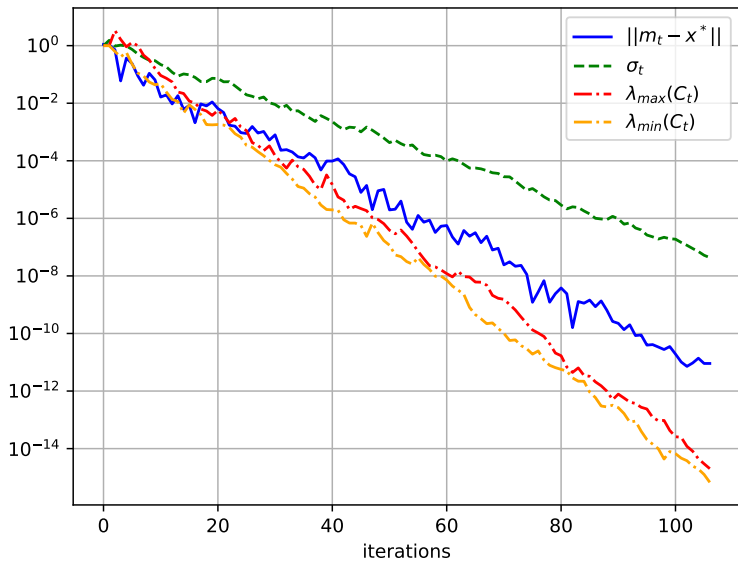
1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$
4. $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$
5. $C_{t+1} = \text{Positive combination}(C_t, \text{Average}[\overleftarrow{(x_{t+1}^{i:\lambda} - m_t)}])$

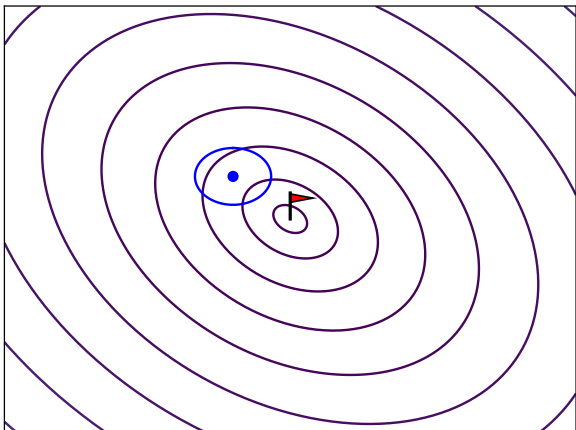


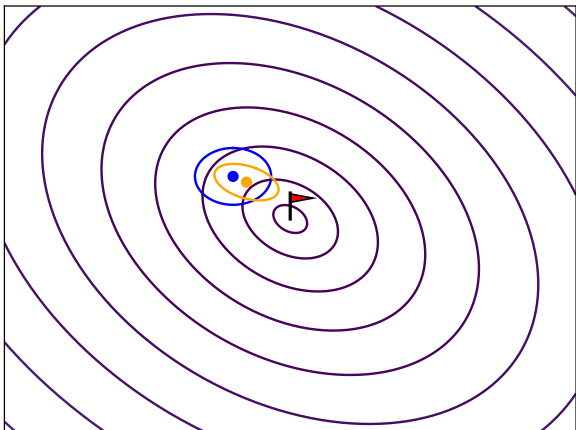
λ = population size

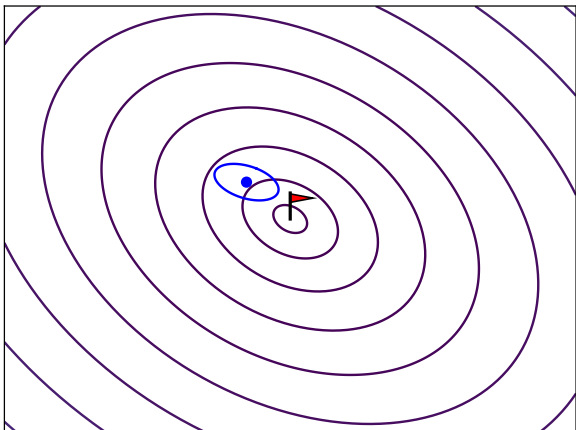
μ = parent number

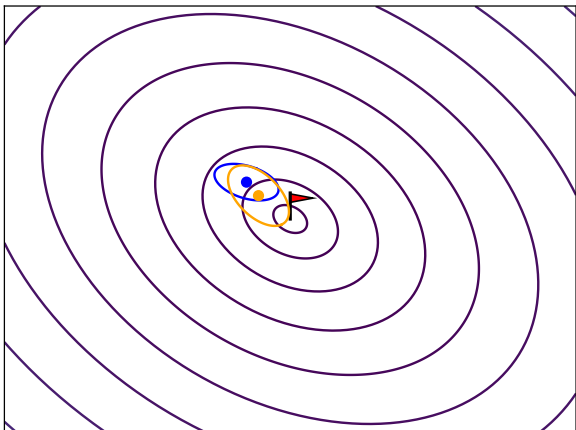
$$f: x \mapsto x^T A x$$

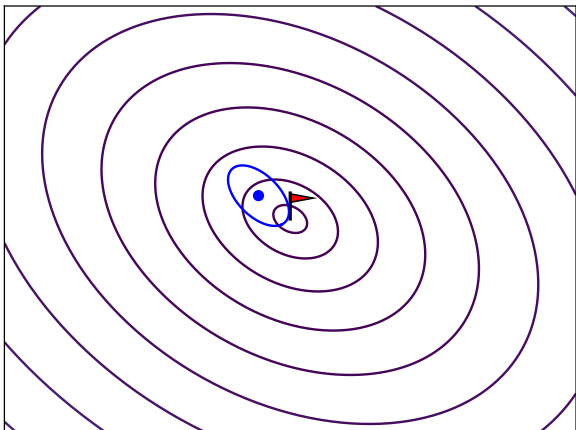


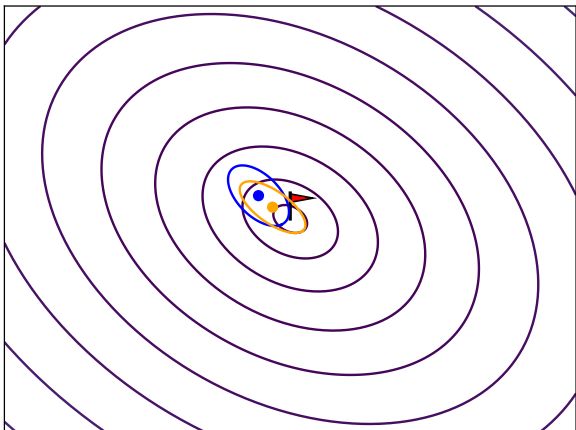


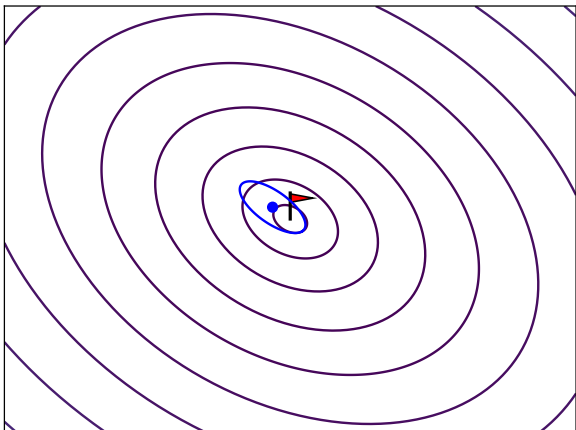


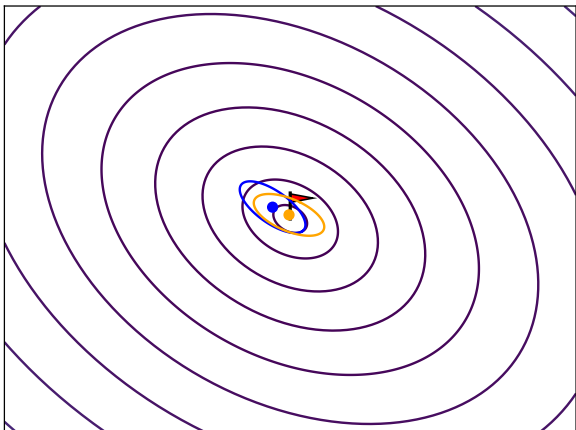


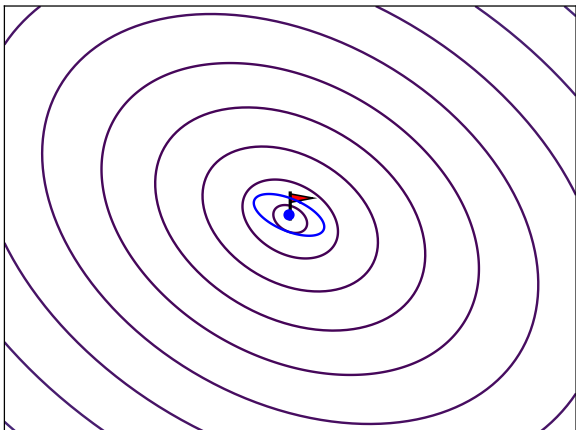


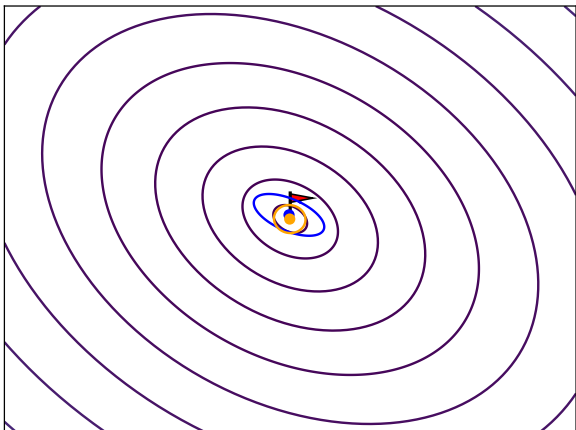


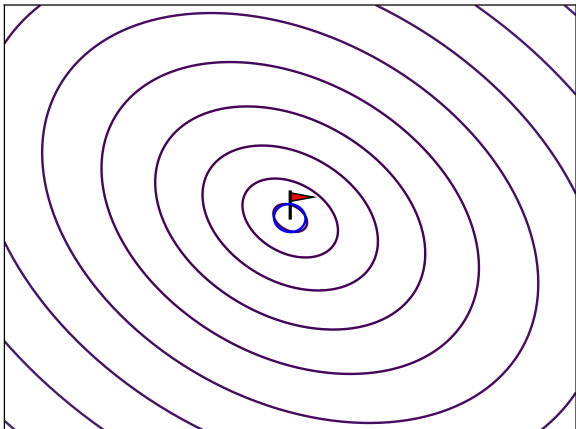


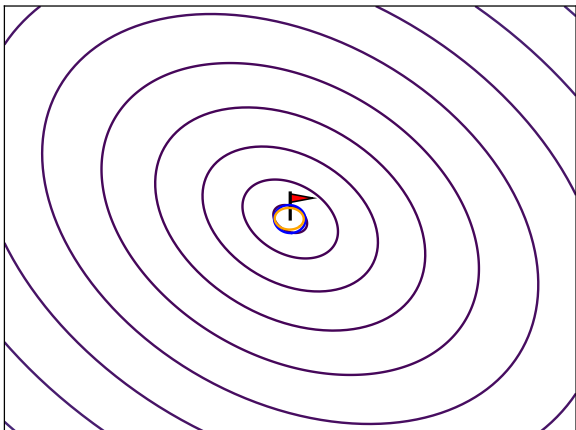


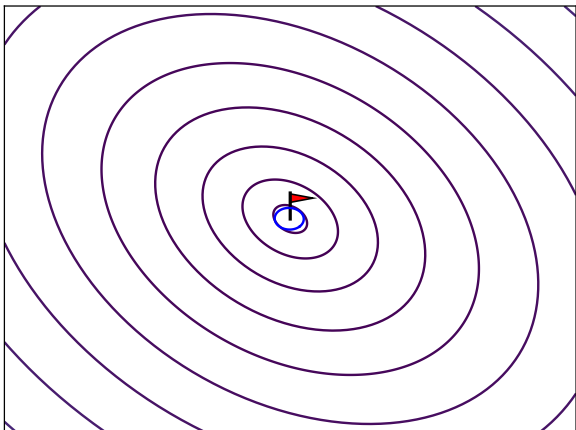


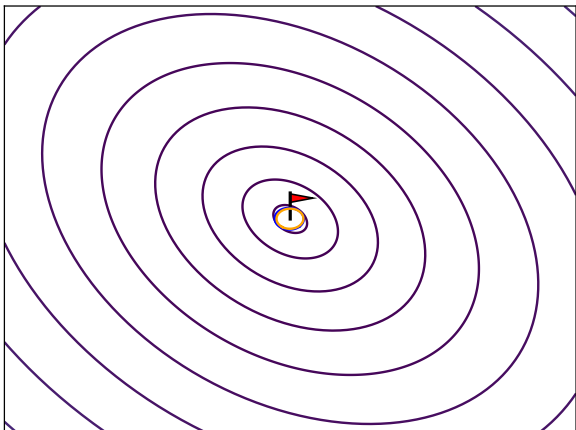


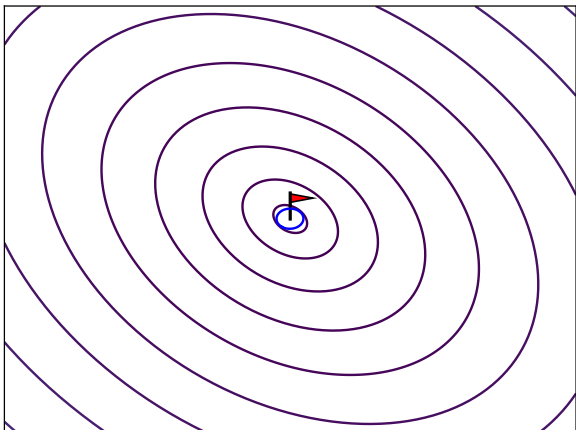


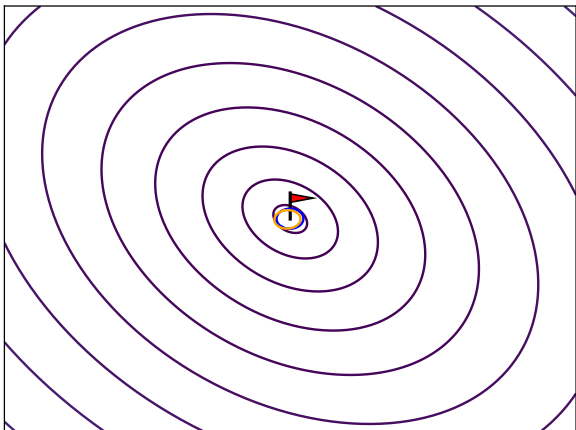


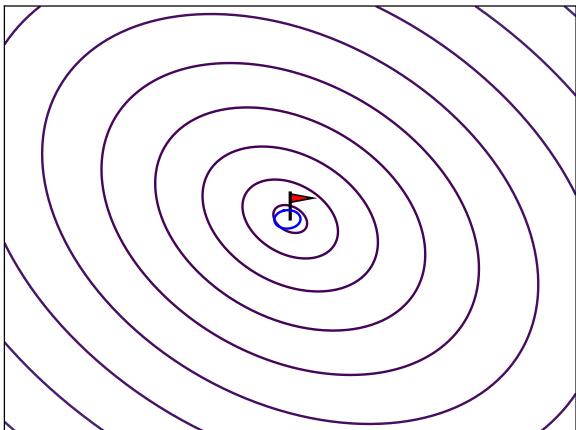


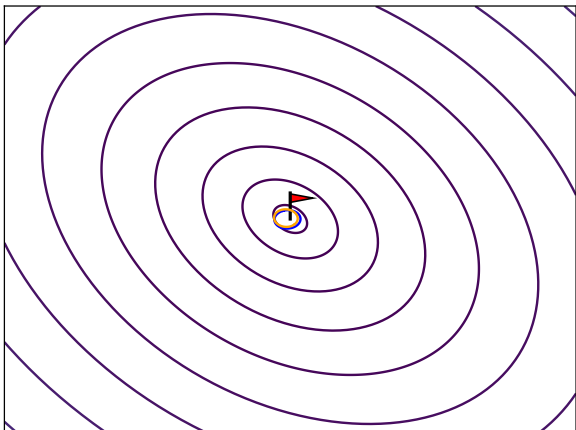


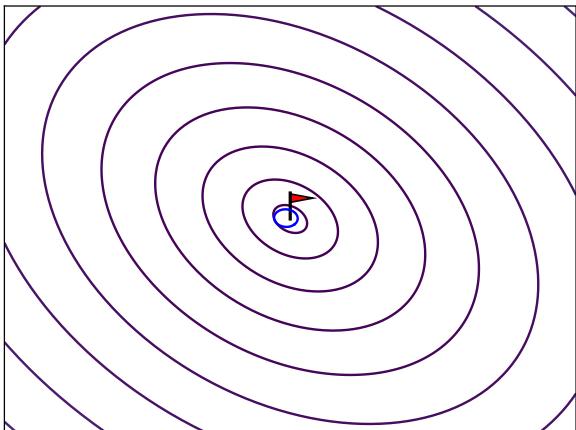


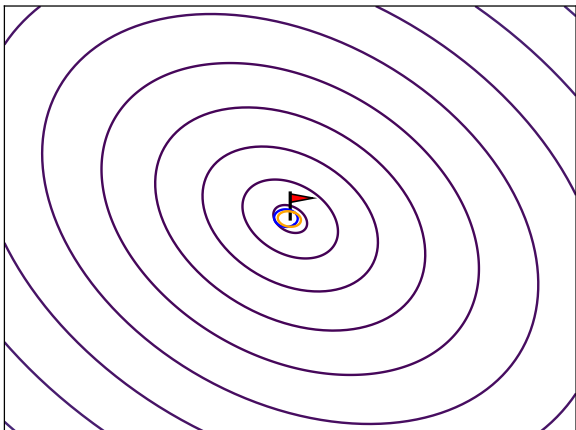


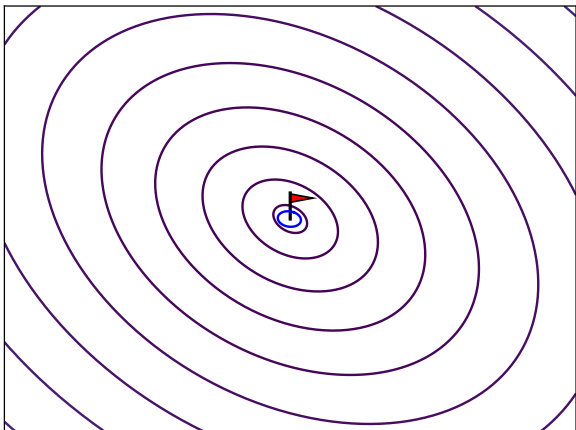


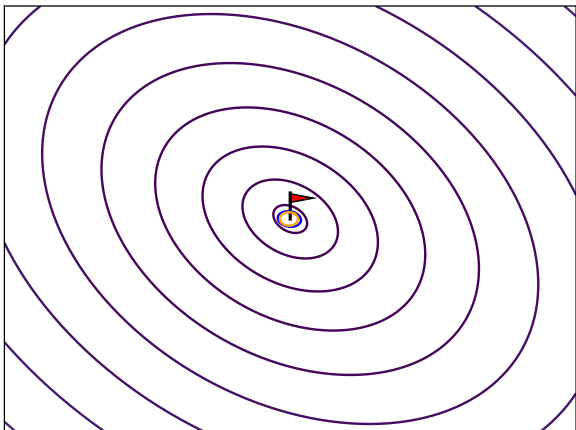


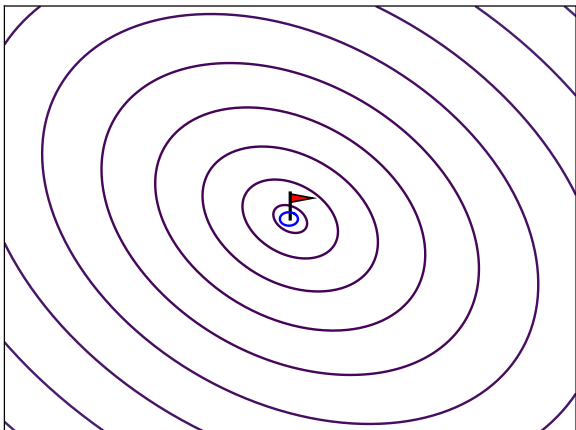


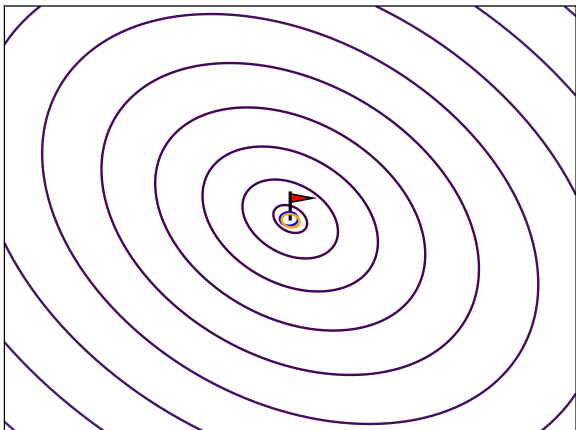


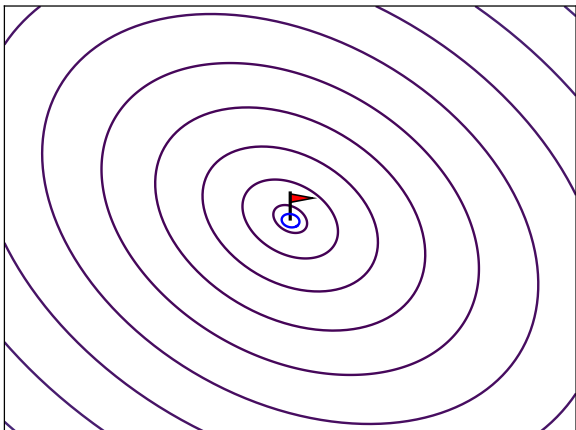


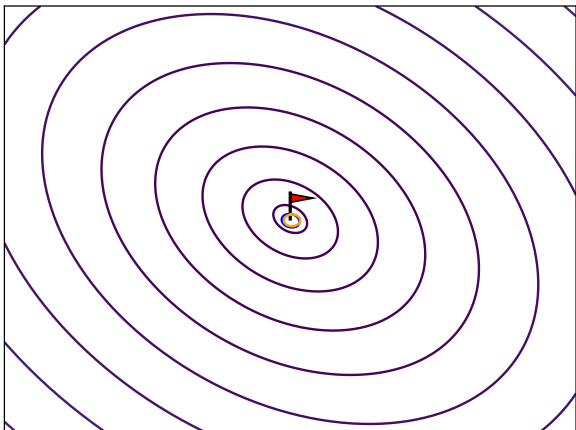












Observations:

$$\lim_{t \rightarrow \infty} m_t = x^*$$

Observations:

$$\lim_{t \rightarrow \infty} m_t = x^*$$

$$\lim_{t \rightarrow \infty} \sigma_t = \lim_{t \rightarrow \infty} C_t = 0$$

Observations:

$$\lim_{t \rightarrow \infty} m_t = x^*$$

$$\lim_{t \rightarrow \infty} \sigma_t = \lim_{t \rightarrow \infty} C_t = 0$$

and

$$\lim_{t \rightarrow \infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}(f)^{-1}$$

Observations:

$$\lim_{t \rightarrow \infty} m_t = x^*$$

$$\lim_{t \rightarrow \infty} \sigma_t = \lim_{t \rightarrow \infty} C_t = 0$$

and

$$\lim_{t \rightarrow \infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}(f)^{-1}$$

For the proof: we rely (again) on Markov chains

$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

$$\Sigma_t = \frac{C_t}{\lambda_{\min}(C_t)}$$

$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

$$\Sigma_t = \frac{C_t}{\lambda_{\min}(C_t)}$$

Proposition

If $f \in \left\{ \begin{array}{c} \text{[contour plot 1]} \\ \text{[contour plot 2]} \\ \text{[contour plot 3]} \\ \text{[contour plot 4]} \end{array} \right\}$, then $\{(z_t, \Sigma_t)\}_{t \in \mathbb{N}}$ is a Markov chain.

Scheme of proof:

1. irreducibility and aperiodicity of $\{(z_t, \Sigma_t)\}_{t \in \mathbb{N}}$
2. drift condition: $\exists K \subset \mathbb{R}^d \times \lambda_{\min}^{-1}(\{1\})$ compact and $V: \mathbb{R}^d \times \lambda_{\min}^{-1}(\{1\}) \rightarrow [1, +\infty]$

$$\mathbb{E}[V(z_1, \Sigma_1)] \leq (1 - \varepsilon)V(z_0, \Sigma_0) \quad \forall (z_0, \Sigma_0) \notin K$$

3. deduce convergence from the ergodicity

Scheme of proof:

1. **irreducibility and aperiodicity of $\{(z_t, \Sigma_t)\}_{t \in \mathbb{N}}$**
2. drift condition: $\exists K \subset \mathbb{R}^d \times \lambda_{\min}^{-1}(\{1\})$ compact and
 $V: \mathbb{R}^d \times \lambda_{\min}^{-1}(\{1\}) \rightarrow [1, +\infty]$

$$\mathbb{E}[V(z_1, \Sigma_1)] \leq (1 - \varepsilon)V(z_0, \Sigma_0) \quad \forall (z_0, \Sigma_0) \notin K$$

3. deduce convergence from the ergodicity

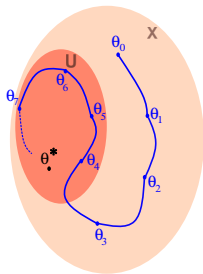
θ^* is **attracting** when

$$\exists x_1, x_2, \dots, \lim_{k \rightarrow \infty} F_k(\theta_0, x_{1..k}) = \theta^*$$

Theorem
If

- $\exists \theta^*$ *attracting*
- $\exists x_1^*, \dots, x_k^*$

such that $F_k(\theta^, \cdot)$ is a **submersion** at $x_{1..k}^*$,*
*then, $\{\theta_t\}_{t \in \mathbb{N}}$ is **irreducible and aperiodic**.*



$$(z_{k+1}, \Sigma_{k+1}) = F(z_k, \Sigma_k, z_{k+1}^{1:\lambda}, \dots, z_{k+1}^{\lambda:\lambda})$$

$$(z_{k+1}, \Sigma_{k+1}) = F(z_k, \Sigma_k, z_{k+1}^{1:\lambda}, \dots, z_{k+1}^{\lambda:\lambda})$$

Proposition

$(0, I_d)$ is attracting

$$(z_{k+1}, \Sigma_{k+1}) = F(z_k, \Sigma_k, z_{k+1}^{1:\lambda}, \dots, z_{k+1}^{\lambda:\lambda})$$

Proposition

$(0, I_d)$ is attracting, and $F_k(0, I_d, \cdot)$ is submersive somewhere.

$$(z_{k+1}, \Sigma_{k+1}) = F(z_k, \Sigma_k, z_{k+1}^{1:\lambda}, \dots, z_{k+1}^{\lambda:\lambda})$$

Proposition

$(0, I_d)$ is attracting, and $F_k(0, I_d, \cdot)$ is submersive somewhere.

Corollary

If $f \in \left\{ \begin{array}{c} \text{[concentric circles]} \\ \text{[elliptical contours]} \\ \text{[wavy lines]} \\ \text{[chaotic pattern]} \end{array} \right\}$, $\{z_t, \Sigma_t\}_{t \in \mathbb{N}}$ is irreducible and aperiodic.

Scheme of proof:

1. irreducibility and aperiodicity of $\{(z_t, \Sigma_t)\}_{t \in \mathbb{N}}$
2. **drift condition:** $\exists K \subset \mathbb{R}^d \times \lambda_{\min}^{-1}(\{1\})$ compact and $V: \mathbb{R}^d \times \lambda_{\min}^{-1}(\{1\}) \rightarrow [1, +\infty]$

$$\mathbb{E}[V(z_1, \Sigma_1)] \leq (1 - \varepsilon)V(z_0, \Sigma_0) \quad \forall (z_0, \Sigma_0) \notin K$$

3. deduce convergence from the ergodicity

$$V(z, \Sigma) = \text{linear combination}(\|z\|^2, \|\Sigma\|)$$

$$V(z, \Sigma) = \text{linear combination}(\|z\|^2, \|\Sigma\|)$$

Proposition

If $f = \text{target}$, then:

$$\mathbb{E}[V(z_{t+1}, \Sigma_{t+1}) \mid z_t, \Sigma_t] \leq (1 - \epsilon) \times V(z_t, \Sigma_t)$$

when $\|z_t\| \gg 1$ or $\|\Sigma_t\| \gg 1$

When $\|\Sigma_t\| \gg \|z_t\|^2$:

When $\lambda_{\max}(\Sigma_t) \gg \|z_t\|^2$:

When $\lambda_{\max}(\Sigma_t) \gg \|z_t\|^2$: we want

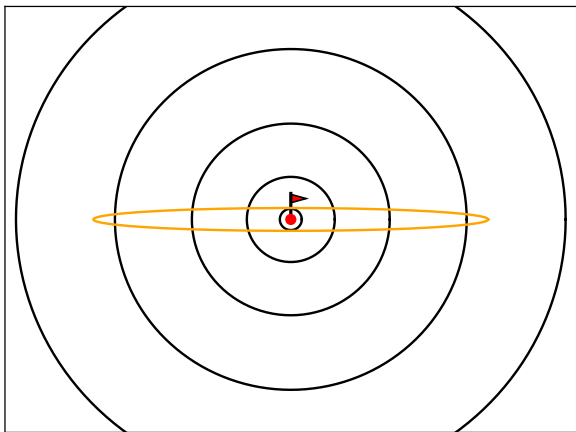
$$\begin{aligned} \mathbb{E}_t [\text{linear combination}(\|z_{t+1}\|^2, \|\Sigma_{t+1}\|)] \\ \leq (1 - \varepsilon) \times \text{linear combination}(\|z_t\|^2, \|\Sigma_t\|) \end{aligned}$$

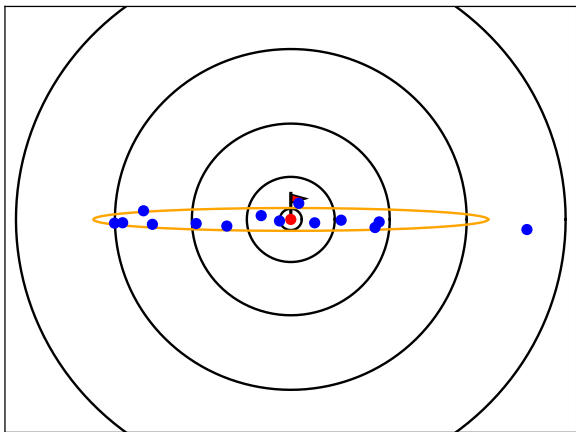
When $\lambda_{\max}(\Sigma_t) \gg \|z_t\|^2$: we want

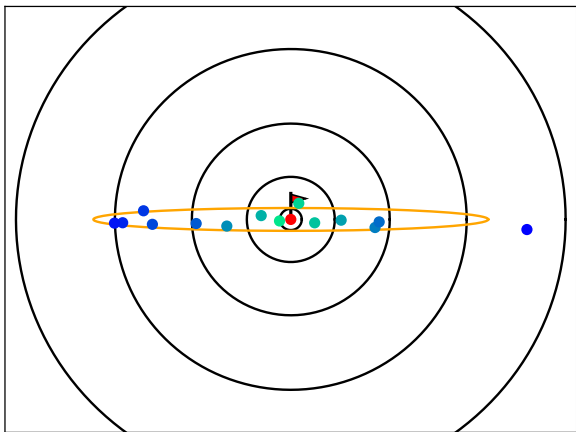
$$\mathbb{E}_t [\|\Sigma_{t+1}\|] \leq (1 - \varepsilon) \times \|\Sigma_t\|$$

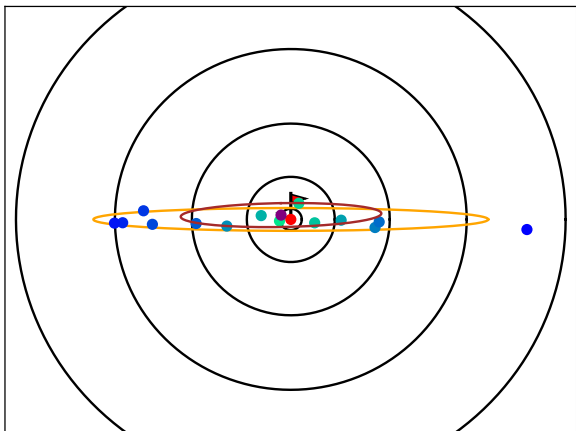
When $\lambda_{\max}(\Sigma_t) \gg \|z_t\|^2$: we want

$$\mathbb{E}_t [\lambda_{\max}(\Sigma_{t+1})] \leq (1 - \varepsilon) \times \lambda_{\max}(\Sigma_t)$$










Proposition

When :

$$\mathbb{E} [\lambda_{\max}(\Sigma_{t+1})] \leq (1 - \varepsilon) \times \lambda_{\max}(\Sigma_t)$$

When $\|z_t\|^2 \gg \|\Sigma_t\|$:

When $\|z_t\|^2 \gg \|\Sigma_t\|$: we want

$$\begin{aligned} \mathbb{E}_t [\text{linear combination}(\|z_{t+1}\|^2, \|\Sigma_{t+1}\|)] \\ \leq (1 - \varepsilon) \times \text{linear combination}(\|z_t\|^2, \|\Sigma_t\|) \end{aligned}$$

When $\|z_t\|^2 \gg \|\Sigma_t\|$: we want

$$\mathbb{E}_t [\|z_{t+1}\|^2] \leq (1 - \varepsilon) \times \|z_t\|^2$$

When $\|z_t\|^2 \gg \|\Sigma_t\|$: we want

$$\mathbb{E}_t [\|z_{t+1}\|^2] \leq (1 - \varepsilon) \times \|z_t\|^2$$

$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

When $\|z_t\|^2 \gg \|\Sigma_t\|$: we want

$$\mathbb{E}_t [\|z_{t+1}\|^2] \leq (1 - \varepsilon) \times \|z_t\|^2$$

$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

$$z_{t+1} = \frac{\text{Update mean}(z_t)}{\text{normalization}}$$

When $\|z_t\|^2 \gg \|\Sigma_t\|$: we want

$$\mathbb{E}_t [\|z_{t+1}\|^2] \leq (1 - \varepsilon) \times \|z_t\|^2$$

$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

$$z_{t+1} = \frac{\text{Update mean}(z_t)}{\text{normalization}}$$

$$\text{normalization} = \frac{\sigma_{t+1}}{\sigma_t} \sqrt{\frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}}$$

When $\|z_t\|^2 \gg \|\Sigma_t\|$: we want

$$\mathbb{E}_t [\|z_{t+1}\|^2] \leq (1 - \varepsilon) \times \|z_t\|^2$$

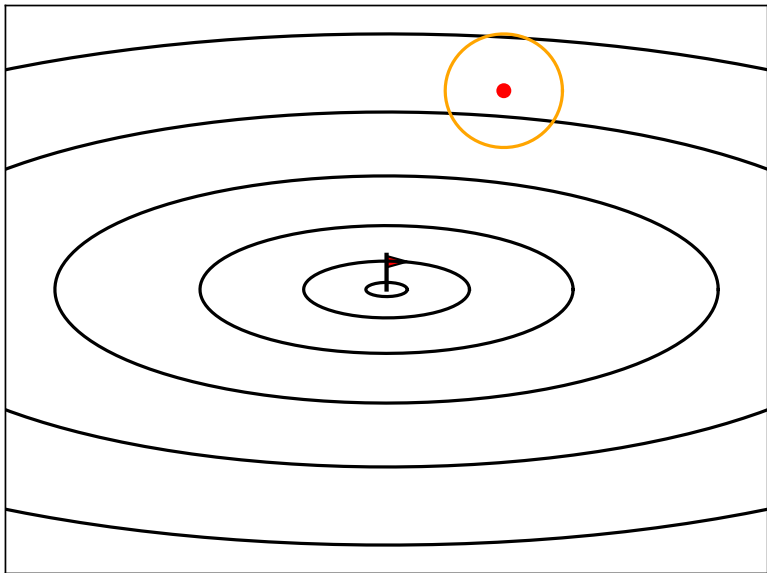
$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

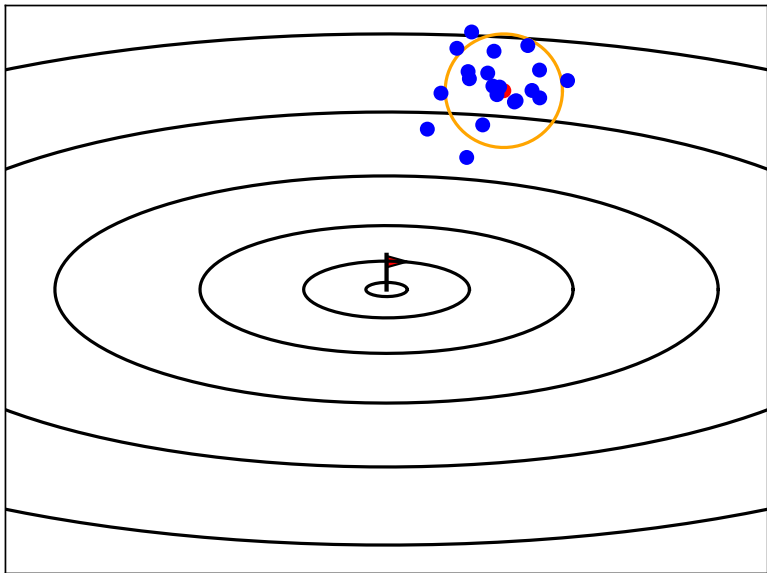
$$z_{t+1} = \frac{\text{Update mean}(z_t)}{\text{normalization}}$$

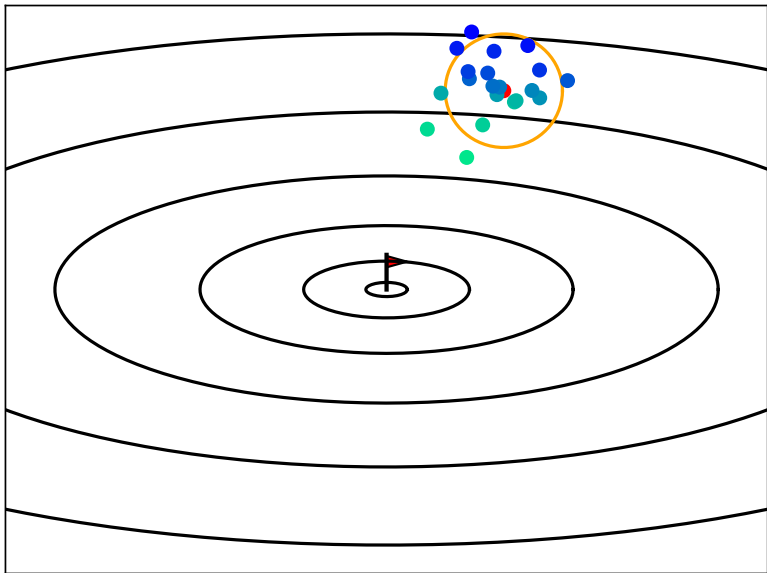
$$\text{normalization} = \underbrace{\frac{\sigma_{t+1}}{\sigma_t}} \sqrt{\frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}}$$

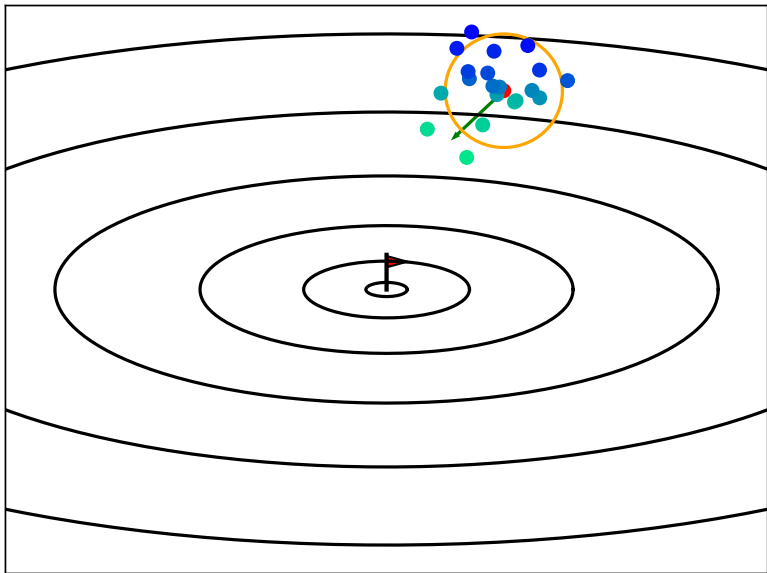
$$\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$$

$$\frac{\sigma_{t+1}}{\sigma_t} = \text{increasing function}(\|z_{t+1} - z_t\|)$$



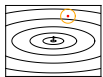






$$\frac{\sigma_{t+1}}{\sigma_t} = \text{increasing function}(\|z_{t+1} - z_t\|)$$

Proposition



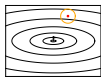
When

:

$$\mathbb{E}[\|z_{t+1} - z_t\|] \geq \text{constant}.$$

$$\frac{\sigma_{t+1}}{\sigma_t} = \text{increasing function}(\|z_{t+1} - z_t\|)$$

Proposition



When

$$\mathbb{E}[\|z_{t+1} - z_t\|] \geq \text{constant}.$$

Corollary

\exists increasing function s.t.:

$$\mathbb{E}[\|z_{t+1}\|^2] \leq (1 - \varepsilon) \times \|z_t\|^2$$

$$V(z, \Sigma) = \text{linear combination}(\|z\|^2, \|\Sigma\|)$$

Proposition

If $f = \text{img}$, then:

$$\mathbb{E}[V(z_{t+1}, \Sigma_{t+1}) \mid z_t, \Sigma_t] \leq (1 - \varepsilon) \times V(z_t, \Sigma_t)$$

when $\|z_t\| \gg 1$ or $\|\Sigma_t\| \gg 1$

Corollary

If $f = \square$, $\{(z_t, \Sigma_t)\}_{t \in \mathbb{N}}$ is ergodic.

Scheme of proof:

1. irreducibility and aperiodicity of $\{(z_t, \Sigma_t)\}_{t \in \mathbb{N}}$
2. drift condition: $\exists K \subset \mathbb{R}^d \times \lambda_{\min}^{-1}(\{1\})$ compact and $V: \mathbb{R}^d \times \lambda_{\min}^{-1}(\{1\}) \rightarrow [1, +\infty]$

$$\mathbb{E}[V(z_1, \Sigma_1)] \leq (1 - \varepsilon)V(z_0, \Sigma_0) \quad \forall (z_0, \Sigma_0) \notin K$$

3. **deduce convergence from the ergodicity**

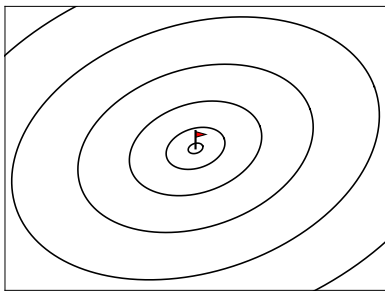
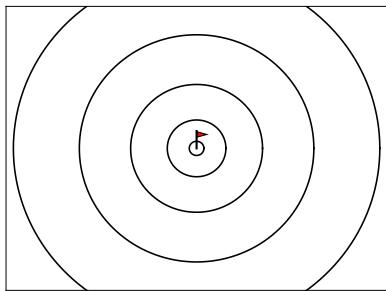
Theorem

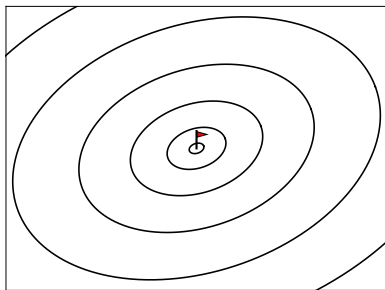
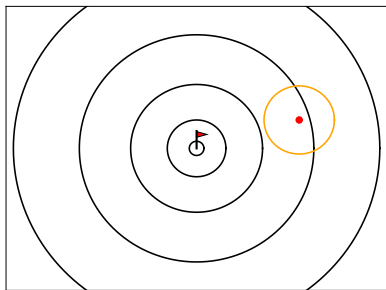
If $f = \text{img}$, CMA-ES converges linearly (or geometrically).

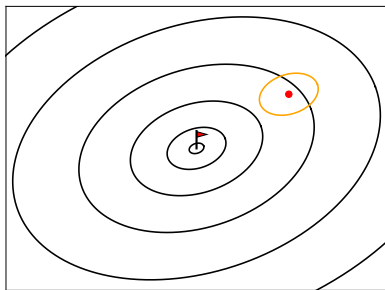
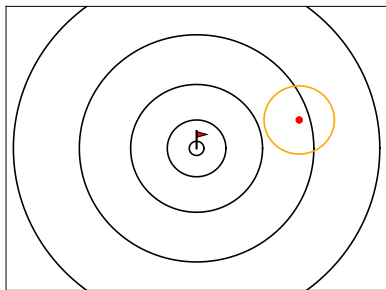
Theorem

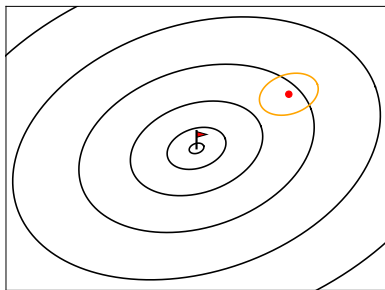
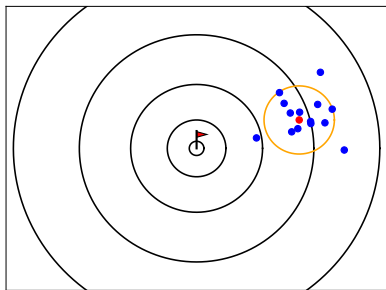
If $f =$ , CMA-ES converges linearly (or geometrically).

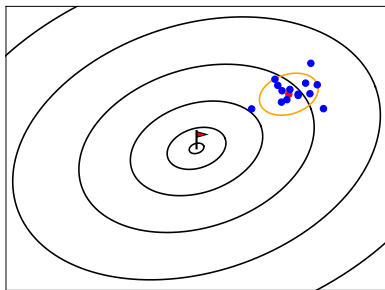
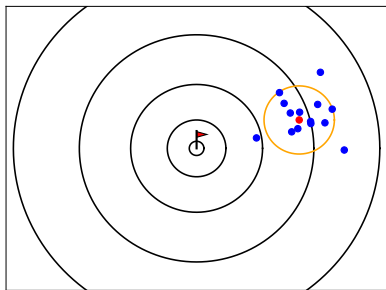
Question: how to extend to $f =$ ?

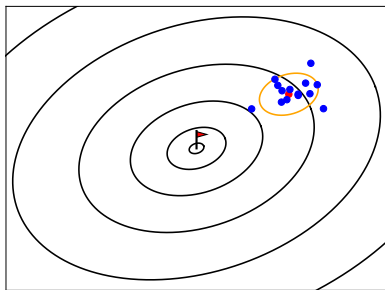
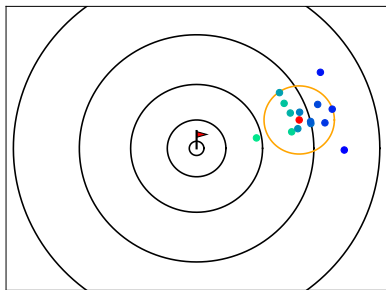


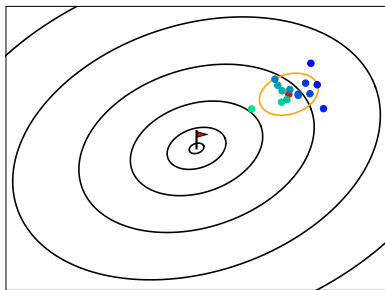
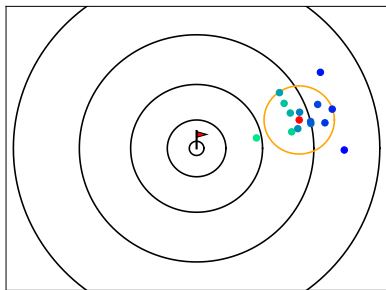


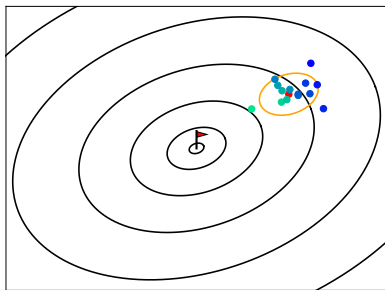
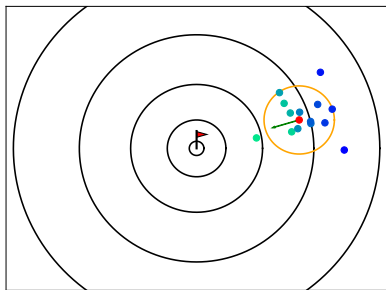


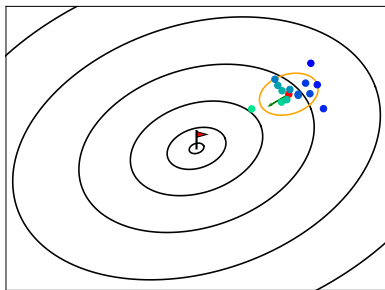
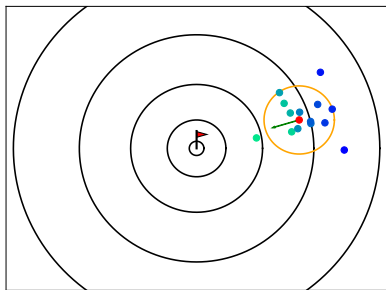


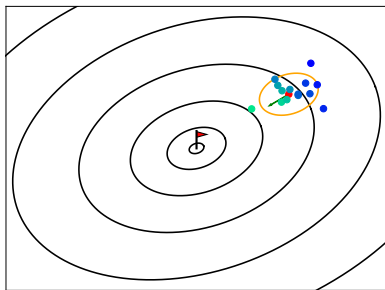
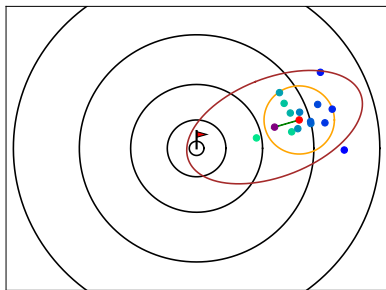


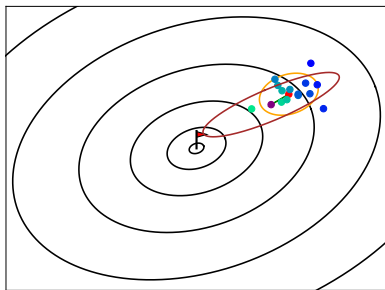
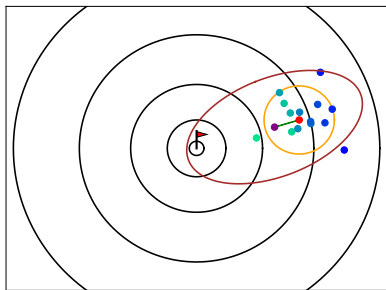


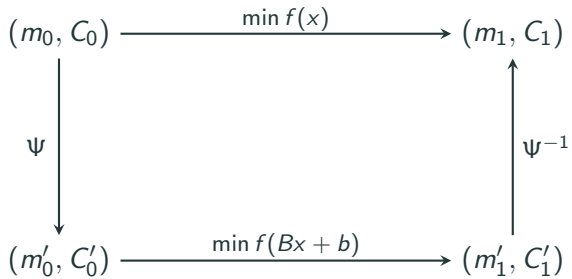










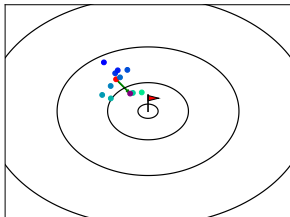


Algorithm 1 Our first ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat: (Given $m_t \in \mathbb{R}^d$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, I_d)$
2. Rank population:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update mean: $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$



λ = population size

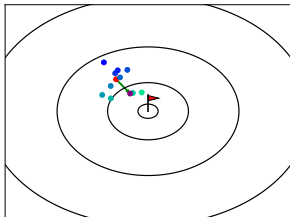
μ = parent number

Algorithm 1 Our first ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat: (Given $m_t \in \mathbb{R}^d$, $C_t \in \mathcal{S}_{++}^d$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, C_t)$
2. Rank population:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update mean: $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$
4. $C_{t+1} = C_t$



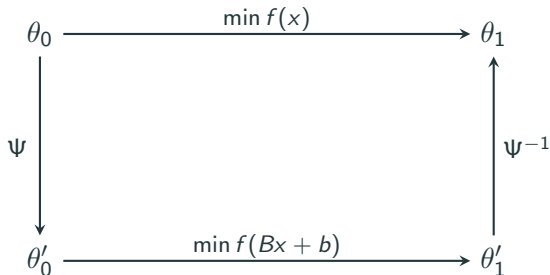
λ = population size

μ = parent number

Definition

An algorithm $\Theta = \{\theta_t\}_{t \in \mathbb{N}}$ is affine-invariant if

(i)



(ii) From θ_0 , Θ can reach a trajectory which starts at $\theta'_0 = \Psi(\theta_0)$.

Theorem

CMA-ES is affine-invariant.

Theorem

CMA-ES is affine-invariant.

Theorem

If $f \in \left\{ \left[\begin{array}{c} \text{circle} \\ \text{ellipse} \end{array} \right] \right\}$:

$$\lim_{t \rightarrow \infty} m_t = x^* \quad \text{geometrically}$$

Theorem

$$f \in \left\{ \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} \right\}:$$

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

Proof.



Theorem

$$f \in \left\{ \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} \right\}:$$

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

Proof.

When $f = \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array}$:



Theorem

$$f \in \left\{ \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} \right\}:$$

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

Proof.

When $f = \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array}$:

$$R \times \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} = \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array}$$

for R a rotation matrix.



Theorem

$$f \in \left\{ \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} \right\}:$$

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

Proof.

When $f = \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array}$:

$$R \times \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} = \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array}$$

for R a rotation matrix.

$$RC_tR^\top \text{ behaves like } C_t$$



Theorem

$$f \in \left\{ \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} \right\}:$$

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

Proof.

When $f = \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array}$:

$$R \times \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} = \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array}$$

for R a rotation matrix.

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{RC_tR^T}{\text{normalization}} \right] = \lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right]$$

□

Theorem

$$f \in \left\{ \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} \right\}:$$

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

Proof.

When $f = \begin{array}{c} \text{[circular contours]} \\ \text{[circular contours]} \end{array}$:

$$R \times \begin{array}{c} \text{[circular contours]} \\ \text{[circular contours]} \end{array} = \begin{array}{c} \text{[circular contours]} \\ \text{[circular contours]} \end{array}$$

for R a rotation matrix.

$$R \times \lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] = \lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \times R$$

□

Theorem

$$f \in \left\{ \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} \right\}:$$

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

Proof.

When $f = \begin{array}{c} \text{[circular contours]} \\ \text{[circular contours]} \end{array}$:

$$R \times \begin{array}{c} \text{[circular contours]} \\ \text{[circular contours]} \end{array} = \begin{array}{c} \text{[circular contours]} \\ \text{[circular contours]} \end{array}$$

for R a rotation matrix.

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto I_d$$

□

Theorem

$$f \in \left\{ \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} \right\}:$$

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

Proof.

When $f = \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array}$:

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto I_d$$

When $f = \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array}$:

□

Theorem

$$f \in \left\{ \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} \right\}:$$

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

Proof.

When $f = \begin{array}{c} \text{[circular contours]} \\ \text{[circular contours]} \end{array}$:

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto I_d$$

When $f = \begin{array}{c} \text{[elliptical contours]} \\ \text{[elliptical contours]} \end{array}$:

$$C_t \text{ behaves like } \text{Hessian}^{-1/2} C_t(\begin{array}{c} \text{[circular contours]} \\ \text{[circular contours]} \end{array}) \text{Hessian}^{-1/2}$$

□

Theorem

$$f \in \left\{ \begin{array}{c} \text{[circular contours]} \\ \text{[elliptical contours]} \end{array} \right\}:$$

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

Proof.

When $f = \begin{array}{c} \text{[circular contours]} \\ \text{[circular contours]} \end{array}$:

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto I_d$$

When $f = \begin{array}{c} \text{[elliptical contours]} \\ \text{[elliptical contours]} \end{array}$:

$$\lim_{t \rightarrow +\infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] \propto \text{Hessian}^{-1}(f)$$

□

$$\sigma_{t+1} = \sigma_t \times \text{increasing function} (\|m_{t+1} - m_t\|)$$

$$\sigma_{t+1} = \sigma_t \times \text{increasing function} (\|m_{t+1} - m_t\|)$$

Goal: remember previous iterations to update σ .

$$\sigma_{t+1} = \sigma_t \times \text{increasing function} (\|m_{t+1} - m_t\|)$$

Goal: remember previous iterations to update σ .

$$\text{path}_{t+1}^\sigma = \text{linear function} (\text{path}_t^\sigma, m_{t+1} - m_t)$$

$$\sigma_{t+1} = \sigma_t \times \text{increasing function} (\|m_{t+1} - m_t\|)$$

Goal: remember previous iterations to update σ .

$$\text{path}_{t+1}^\sigma = \text{linear function} (\text{path}_t^\sigma, m_{t+1} - m_t)$$

New update:

$$\sigma_{t+1} = \sigma_t \times \text{increasing function} (\|\text{path}_{t+1}^\sigma\|)$$

$$\text{path}_{t+1}^c = \text{linear function}(\text{path}_t^c, m_{t+1} - m_t)$$

$$\text{path}_{t+1}^c = \text{linear function}(\text{path}_t^c, m_{t+1} - m_t)$$

$$C_{t+1} = \text{Linear combination} \left(C_t, \text{Average}[\overleftarrow{x_{t+1}^{i:\lambda}} - m_t], \overleftarrow{\text{path}_{t+1}^c} \right)$$

Algorithm 5 CMA-ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

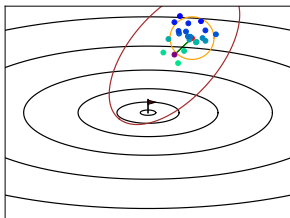
Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$
- 4.

$$\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$$

5.

$$C_{t+1} = \text{Linear combination} \left(C_t, \text{Average} \left[\overleftarrow{(x_{t+1}^{i:\lambda} - m_t)} \right] \right)$$



λ = population size

μ = parent number

Algorithm 5 CMA-ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$

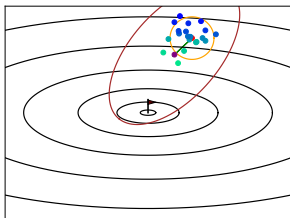
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$

3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$

4. $\text{path}_{t+1}^\sigma = \text{Linear}(\text{path}_t^\sigma, m_{t+1} - m_t)$
 $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|m_{t+1} - m_t\|)$

5.

$$C_{t+1} = \text{Linear combination} \left(C_t, \text{Average} \left[\overleftarrow{x_{t+1}^{i:\lambda} - m_t} \right] \right)$$



λ = population size

μ = parent number

Algorithm 5 CMA-ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$

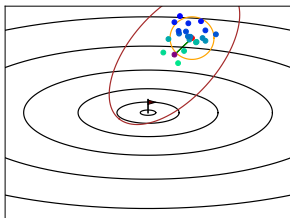
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$

3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$

4. $\text{path}_{t+1}^\sigma = \text{Linear}(\text{path}_t^\sigma, m_{t+1} - m_t)$
 $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|\text{path}_{t+1}^\sigma\|)$

5.

$$C_{t+1} = \text{Linear combination} \left(C_t, \text{Average} \left[\overleftarrow{(x_{t+1}^{i:\lambda} - m_t)} \right] \right)$$



λ = population size

μ = parent number

Algorithm 5 CMA-ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$

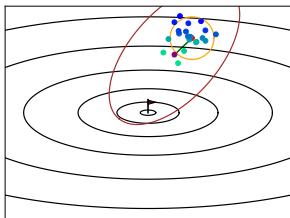
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$

3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$

4. $\text{path}_{t+1}^\sigma = \text{Linear}(\text{path}_t^\sigma, m_{t+1} - m_t)$
 $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|\text{path}_{t+1}^\sigma\|)$

5. $\text{path}_{t+1}^c = \text{Linear}(\text{path}_t^c, m_{t+1} - m_t)$

$C_{t+1} = \text{Linear combination} \left(C_t, \text{Average} \left[\overleftarrow{x_{t+1}^{i:\lambda} - m_t} \right] \right)$



λ = population size

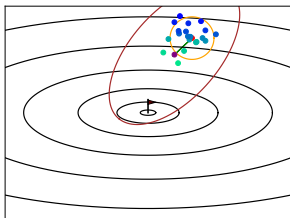
μ = parent number

Algorithm 5 CMA-ES

Goal: $\min_{x \in \mathbb{R}^d} f(x)$

Repeat ($m_t \in \mathbb{R}^d$, $\sigma_t > 0$, $C_t \succ 0$)

1. $x_{t+1}^1, \dots, x_{t+1}^\lambda \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$
2. sort $f(x_{t+1}^i)$:
 $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. $m_{t+1} = \text{Average}(x_{t+1}^{1:\lambda}, \dots, x_{t+1}^{\mu:\lambda})$



4. $\text{path}_{t+1}^\sigma = \text{Linear}(\text{path}_t^\sigma, m_{t+1} - m_t)$
 $\sigma_{t+1} = \sigma_t \times \text{increasing function}(\|\text{path}_{t+1}^\sigma\|)$
5. $\text{path}_{t+1}^c = \text{Linear}(\text{path}_t^c, m_{t+1} - m_t)$
 $C_{t+1} = \text{Linear combination} \left(C_t, \text{Average} \left[\overleftarrow{(x_{t+1}^{i:\lambda} - m_t)} \right], \overleftarrow{\text{path}_{t+1}^c} \right)$

λ = population size

μ = parent number

$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

$$\Sigma_t = \frac{C_t}{\lambda_{\min}(C_t)}$$

$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

$$\Sigma_t = \frac{C_t}{\lambda_{\min}(C_t)}$$

$$\text{normalized path}_t^{c,\sigma} = \frac{\text{path}_t^{c,\sigma}}{\text{normalization}}$$

$$z_t = \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

$$\Sigma_t = \frac{C_t}{\lambda_{\min}(C_t)}$$

$$\text{normalized path}_t^{c,\sigma} = \frac{\text{path}_t^{c,\sigma}}{\text{normalization}}$$

Proposition

If $f \in \left\{ \begin{array}{c} \text{[contour plot]} \\ \text{[contour plot]} \\ \text{[contour plot]} \\ \text{[contour plot]} \end{array} \right\}$, then $\{(z_t, \Sigma_t, n. \text{ path}_t^c, n. \text{ path}_t^\sigma)\}_{t \in \mathbb{N}}$ is a Markov chain.

When X is infinite:

Theorem

If $\{\theta_t\}_{t \in \mathbb{N}}$ is an irreducible, aperiodic Markov chain, then it is ergodic if

$$\mathbb{E}[V(\theta_{t+1}) \mid \theta_t] \leq (1 - \varepsilon)V(\theta_t) \quad \text{if } \theta_t \notin \text{compact set}$$

for some $V: X \rightarrow [1, +\infty]$.

When X is infinite:

Theorem

If $\{\theta_t\}_{t \in \mathbb{N}}$ is an irreducible, aperiodic Markov chain, then it is ergodic if

$$\mathbb{E}[V(\theta_{t+n(\theta_t)}) \mid \theta_t] \leq (1 - \varepsilon)^{n(\theta_t)} V(\theta_t) \quad \text{if } \theta_t \notin \text{compact set}$$

for some $V: X \rightarrow [1, +\infty]$.

$$V(z, \Sigma, \text{n. path}^c, \text{n. path}^\sigma) = \text{linear combination}(\|z\|^2, \|\Sigma\|, \|\text{n. paths}\|^2)$$

$$V(z, \Sigma, n. \text{ path}^c, n. \text{ path}^\sigma) = \text{linear combination}(\|z\|^2, \|\Sigma\|, \|n. \text{ paths}\|^2)$$

Proposition

If $f = \text{img}$, then:

$$\mathbb{E}[V(\text{iteration}_{t+k}) \mid \text{iteration}_t] \leq (1 - \varepsilon) \times V(\text{iteration}_t)$$

when $\|z_t\| \gg 1$ or $\|\Sigma_t\| \gg 1$ or $\|n. \text{ path}_t^c\| \gg 1$ or $\|n. \text{ path}_t^\sigma\| \gg 1$.

Theorem

If $f \in \left\{ \begin{array}{c} \text{[contour plot 1]} \\ \text{[contour plot 2]} \end{array} \right\}$:

$$\lim_{t \rightarrow \infty} m_t = x^* \quad \text{geometrically}$$

and

$$\lim_{t \rightarrow \infty} \mathbb{E} \left[\frac{C_t}{\text{normalization}} \right] = H^{-1}$$

Thank you!