

Convergence Analysis of Evolution Strategies with Covariance Matrix Adaptation (CMA-ES) via Markov Chain Stability Analysis

Blackbox Optimization and Derivative-Free Algorithms

Armand Gissler

Friday 2nd June, 2023

CMA, École polytechnique & Inria
(with Anne Auger & Nikolaus Hansen)



Inria

Consider the optimisation problem

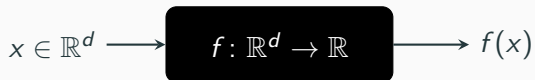
$$\min_{x \in \mathbb{R}^d} f(x) \quad (\text{P})$$

Black-box optimisation and Evolution strategies

Consider the optimisation problem

$$\min_{x \in \mathbb{R}^d} f(x) \quad (\text{P})$$

with

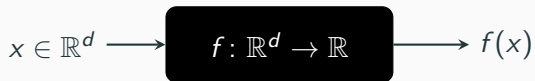


Black-box optimisation and Evolution strategies

Consider the optimisation problem

$$\min_{x \in \mathbb{R}^d} f(x) \quad (\text{P})$$

with



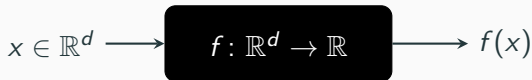
\Rightarrow we only have access to a minimum amount of informations on f

Black-box optimisation and Evolution strategies

Consider the optimisation problem

$$\min_{x \in \mathbb{R}^d} f(x) \quad (\text{P})$$

with



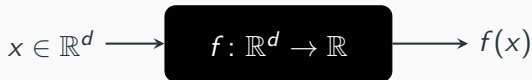
\Rightarrow we only have access to a minimum amount of informations on f
(in particular no information on the derivatives of f)

Black-box optimisation and Evolution strategies

Consider the optimisation problem

$$\min_{x \in \mathbb{R}^d} f(x) \quad (\text{P})$$

with



\Rightarrow we only have access to a minimum amount of informations on f
(in particular no information on the derivatives of f)

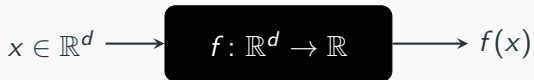
CMA-ES approximates the minimum x^* of f by a multivariate normal distribution $\mathcal{N}(m, \sigma^2 C)$

Black-box optimisation and Evolution strategies

Consider the optimisation problem

$$\min_{x \in \mathbb{R}^d} f(x) \quad (\text{P})$$

with



\Rightarrow we only have access to a minimum amount of informations on f
(in particular no information on the derivatives of f)

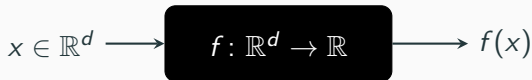
CMA-ES approximates the minimum x^* of f by a multivariate normal distribution $\mathcal{N}(m, \sigma^2 C)$ **by adapting the mean** $m \in \mathbb{R}^d$

Black-box optimisation and Evolution strategies

Consider the optimisation problem

$$\min_{x \in \mathbb{R}^d} f(x) \quad (\text{P})$$

with



\Rightarrow we only have access to a minimum amount of informations on f
(in particular no information on the derivatives of f)

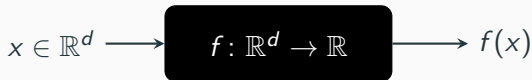
CMA-ES approximates the minimum x^* of f by a multivariate normal distribution $\mathcal{N}(m, \sigma^2 C)$ **by adapting** the mean $m \in \mathbb{R}^d$, **the stepsize** $\sigma > 0$

Black-box optimisation and Evolution strategies

Consider the optimisation problem

$$\min_{x \in \mathbb{R}^d} f(x) \quad (\text{P})$$

with



\Rightarrow we only have access to a minimum amount of informations on f
(in particular no information on the derivatives of f)

CMA-ES approximates the minimum x^* of f by a multivariate normal distribution $\mathcal{N}(m, \sigma^2 C)$ **by adapting** the mean $m \in \mathbb{R}^d$, the stepsize $\sigma > 0$ **and the covariance matrix** $C \in \mathcal{S}_{++}^d$.

- converges linearly to the minimum

- converges linearly to the minimum
- **learns the inverse Hessian** of a convex-quadratic function

- converges linearly to the minimum
- learns the inverse Hessian of a convex-quadratic function
- state-of-the-art method for **difficult (derivative-free) optimization problems** such as minimization of **ill-conditioned, non-separable, discontinuous, multimodal** and/or **noisy** functions

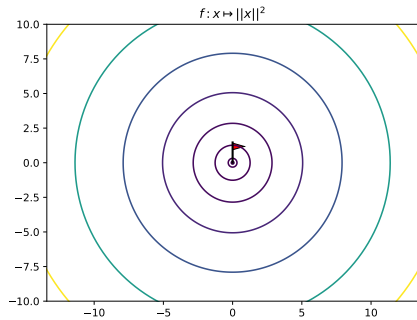
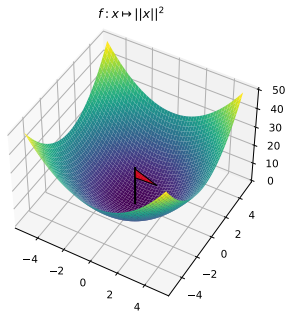
- converges linearly to the minimum
- learns the inverse Hessian of a convex-quadratic function
- state-of-the-art method for difficult (derivative-free) optimization problems such as minimization of ill-conditioned, non-separable, discontinuous, multimodal and/or noisy functions
- **> 40 millions of downloads of two Python modules**

- converges linearly to the minimum
- learns the inverse Hessian of a convex-quadratic function
- state-of-the-art method for difficult (derivative-free) optimization problems such as minimization of ill-conditioned, non-separable, discontinuous, multimodal and/or noisy functions
- > 40 millions of downloads of two Python modules
- **proofs of convergence require additional assumptions so far**

- Without covariance matrix adaptation: Touré et al, *Global linear convergence of Evolution Strategies with recombination on scaling-invariant functions* (2021)
- With a sufficient decrease condition: Diouane et al, *Globally convergent evolution strategies* (2015)
- Assuming that the covariance matrix is bounded: Akimoto et al, *Global linear convergence of evolution strategies on more than smooth strongly convex functions* (2022)
- Using a different update for the covariance matrix: Glasmachers et al, *Convergence analysis of the Hessian estimation evolution strategy* (2022)

CMA-ES: algorithm presentation

Level sets representation



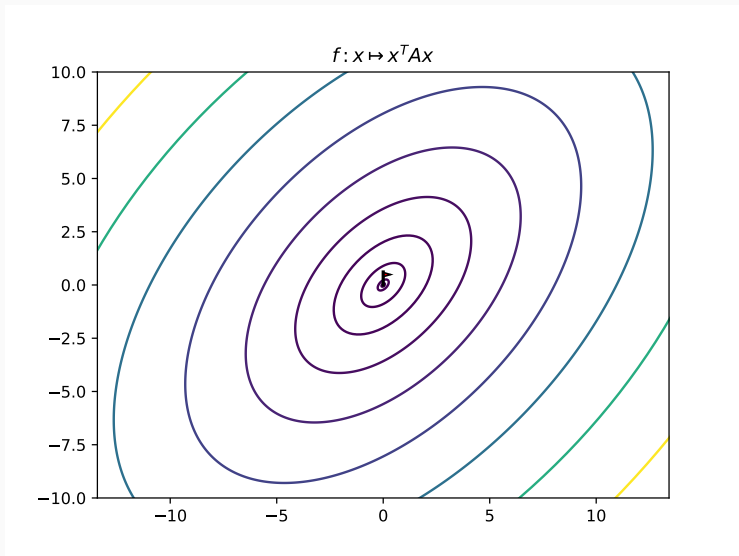
CMA-ES : presentation

Principle of Evolution Strategies (ES) : approximate the minimum of the function by a distribution $\mathcal{N}(m, \sigma^2 C)$.

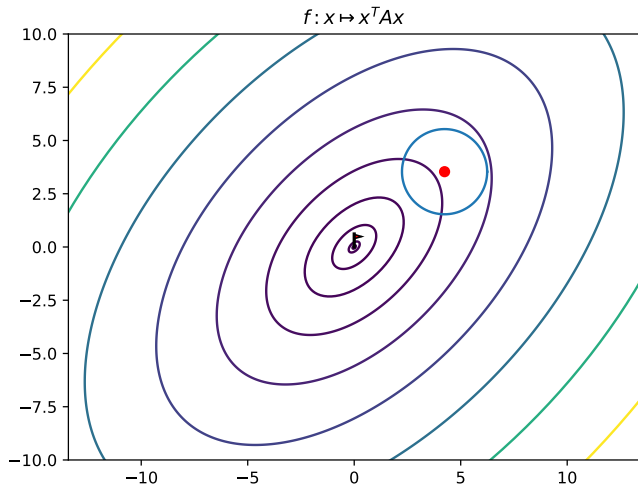
CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

Start from a distribution $\mathcal{N}(m_t, \sigma_t^2 C_t)$



Start from a distribution $\mathcal{N}(m_t, \sigma_t^2 C_t)$



CMA-ES : presentation

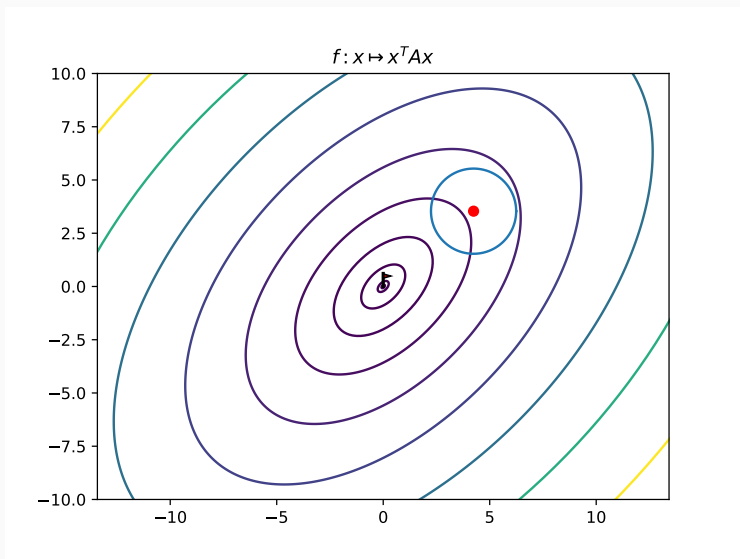
At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

CMA-ES : presentation

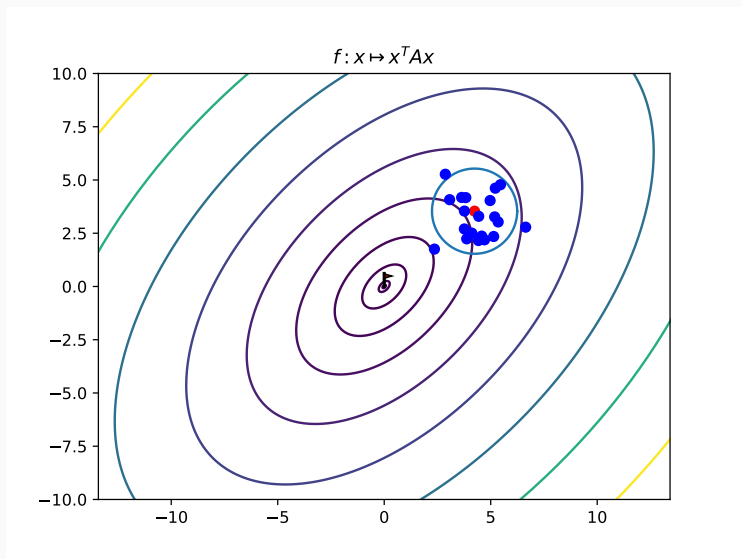
At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. **Generate** λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*

Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ independently



Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ independently



CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. **Generate** λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

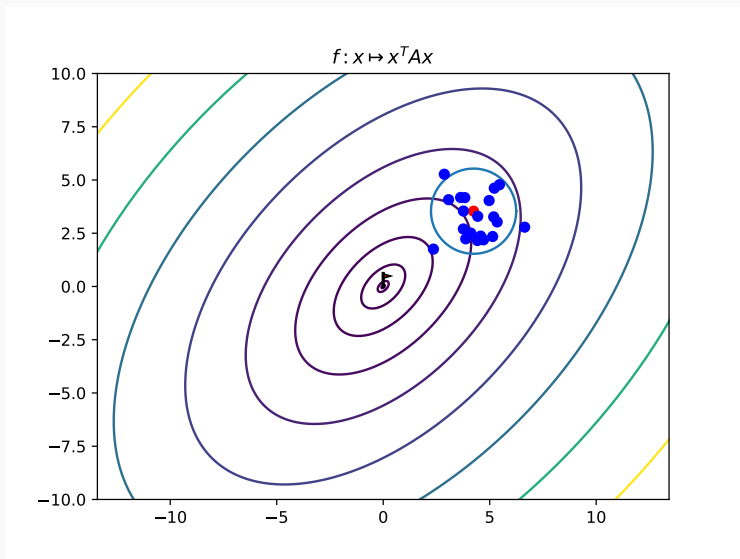
1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. **Rank** the x_{t+1}^i w.r.t. their f -values

CMA-ES : presentation

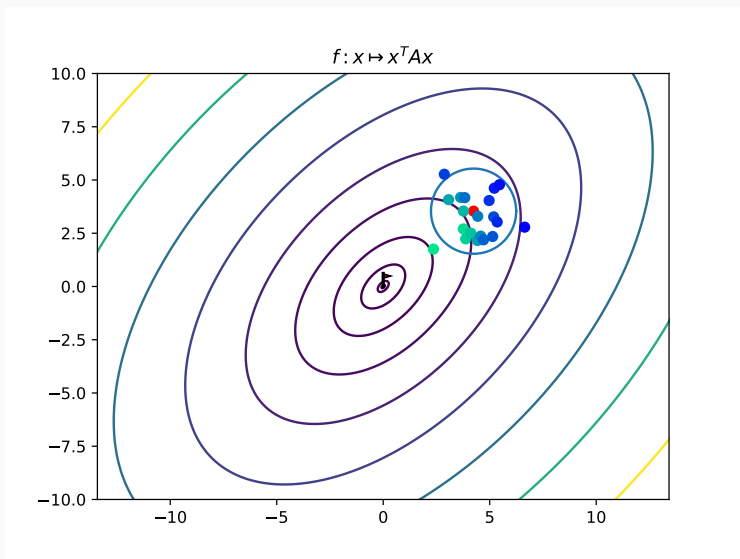
At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. **Rank** the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$

Rank the x_{t+1}^i w.r.t. their f -values $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda})$



Rank the x_{t+1}^i w.r.t. their f -values $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda})$



CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. **Rank** the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. **Update the mean**

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. **Update the mean:** $m_{t+1} = \text{Weighted average}(\text{population})$

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. **Update the mean:** $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. **Update the mean:** $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

The best offspring have the largest weights

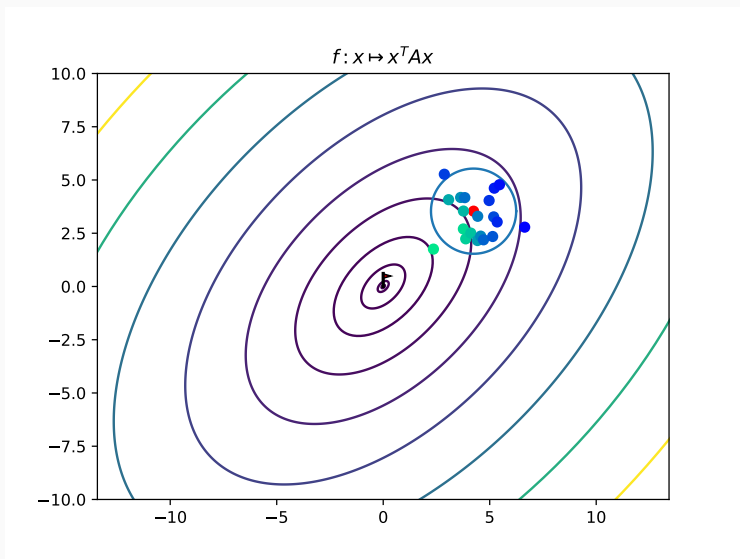
CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

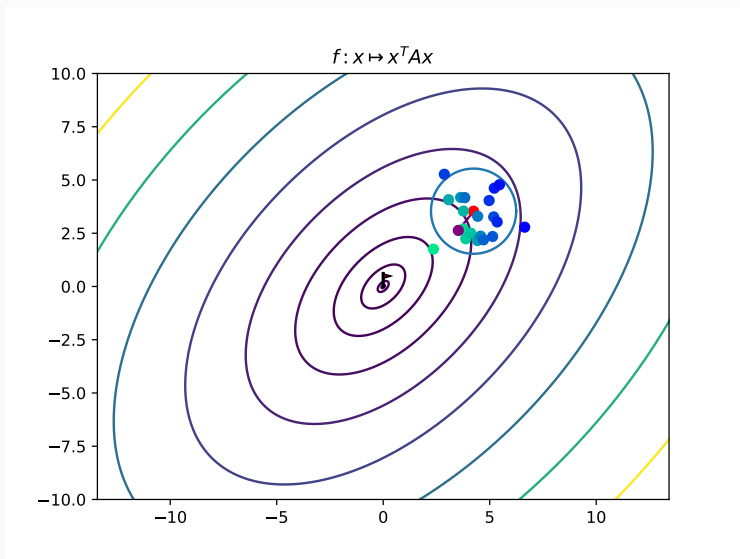
1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. **Update the mean:** $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

The best offspring have the largest weights: $w_1 \geq w_2 \dots$

Update the mean $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$



Update the mean $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$



CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. **Update the mean:** $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

The best offspring have the largest weights: $w_1 \geq w_2 \dots$

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update the mean: $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

The best offspring have the largest weights: $w_1 \geq w_2 \dots$

4. **Update the stepsize**

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

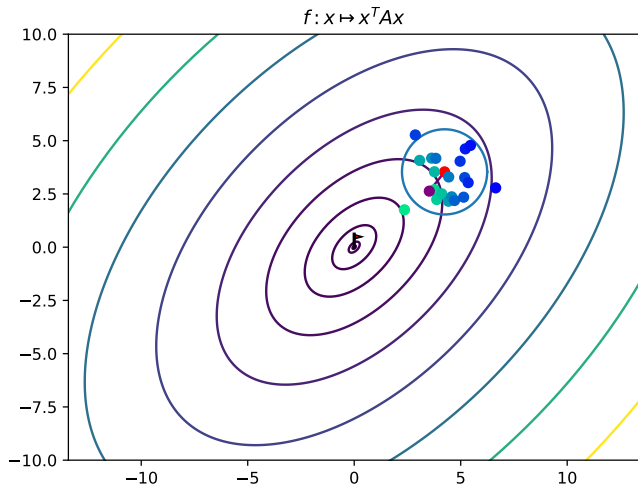
1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update the mean: $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

The best offspring have the largest weights: $w_1 \geq w_2 \dots$

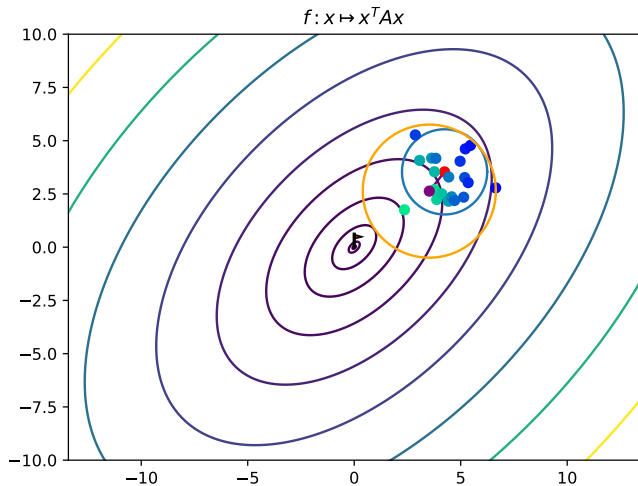
4. **Update the stepsize:**

Increase the stepsize if the path taken by the mean is **larger than expected** (assuming no selection)

Adapt the stepsize



Adapt the stepsize



CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update the mean: $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

The best offspring have the largest weights: $w_1 \geq w_2 \dots$

4. **Update the stepsize:**

Increase the stepsize if the path taken by the mean is **larger than expected** (assuming no selection)

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update the mean: $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

The best offspring have the largest weights: $w_1 \geq w_2 \dots$

4. Update the stepsize:
Increase the stepsize if the path taken by the mean is larger than expected (assuming no selection)
5. **Update the covariance matrix**

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update the mean: $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

The best offspring have the largest weights: $w_1 \geq w_2 \dots$

4. Update the stepsize:
Increase the stepsize if the path taken by the mean is larger than expected (assuming no selection)
5. **Update the covariance matrix**

$$C_{t+1} = (1 - c)C_t + c' \frac{(m_{t+1} - m_t)(m_{t+1} - m_t)^T}{\sigma_t^2}$$

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update the mean: $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

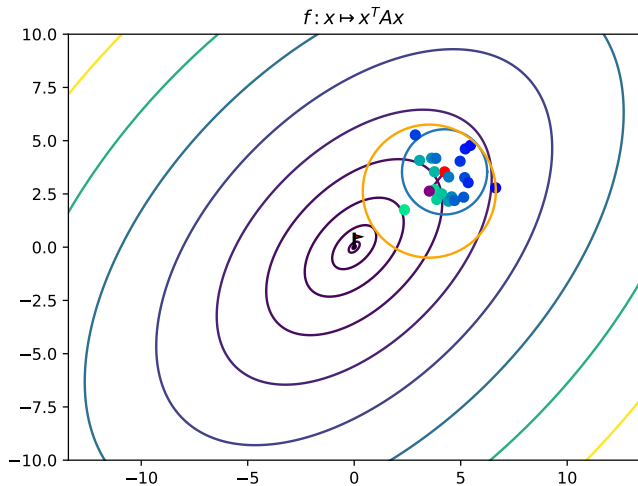
The best offspring have the largest weights: $w_1 \geq w_2 \dots$

4. Update the stepsize:
Increase the stepsize if the path taken by the mean is larger than expected (assuming no selection)

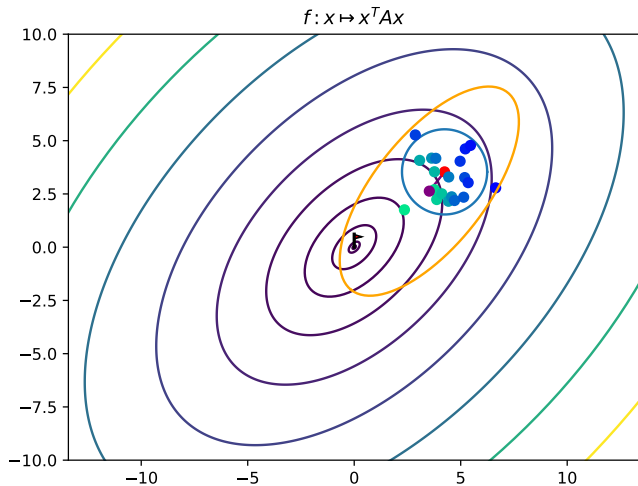
5. **Update the covariance matrix**

$$C_{t+1} = (1 - c)C_t + c \sum w_i \frac{(x_{t+1}^{i:\lambda} - m_t)(x_{t+1}^{i:\lambda} - m_t)^T}{\sigma_t^2}$$

Adapt the covariance matrix



Adapt the covariance matrix



CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

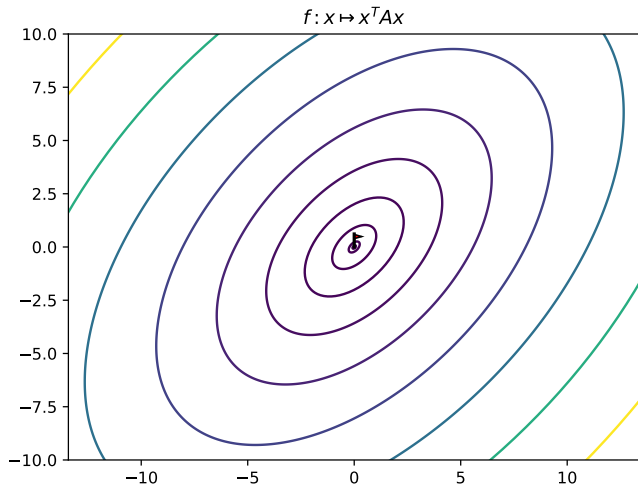
1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update the mean: $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

The best offspring have the largest weights: $w_1 \geq w_2 \dots$

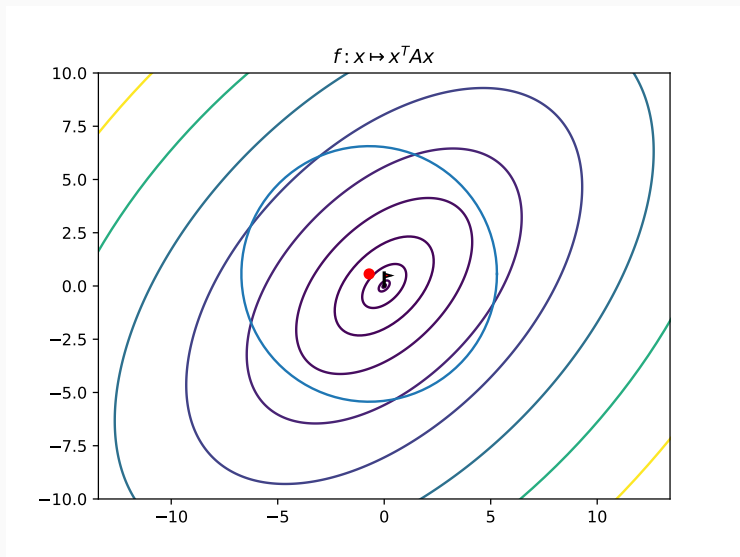
4. Update the stepsize:
Increase the stepsize if the path taken by the mean is larger than expected (assuming no selection)
5. Update the covariance matrix

$$C_{t+1} = (1 - c)C_t + c \sum w_i \frac{(x_{t+1}^{i:\lambda} - m_t)(x_{t+1}^{i:\lambda} - m_t)^T}{\sigma_t^2}$$

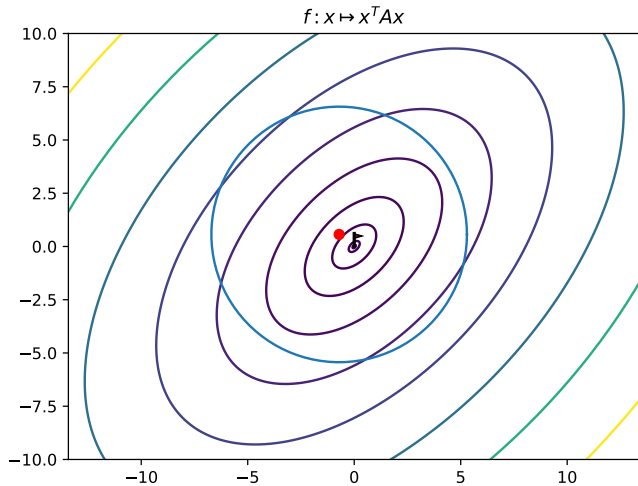
Start from a distribution $\mathcal{N}(m_t, \sigma_t^2 C_t)$



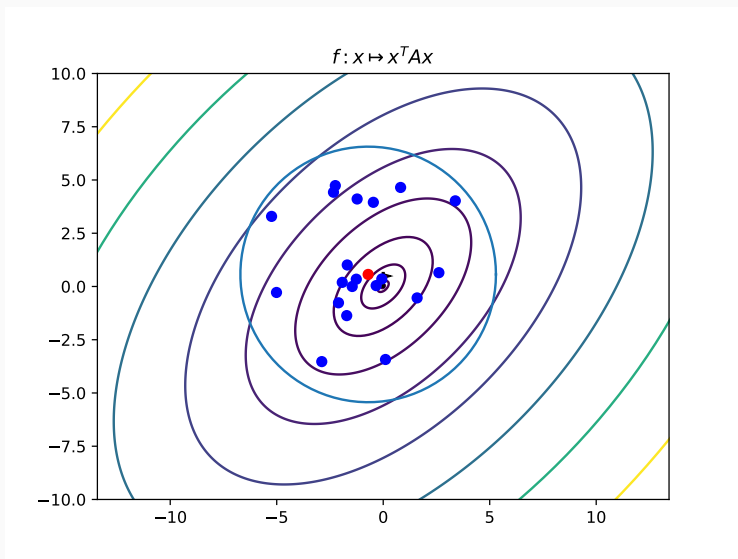
Start from a distribution $\mathcal{N}(m_t, \sigma_t^2 C_t)$



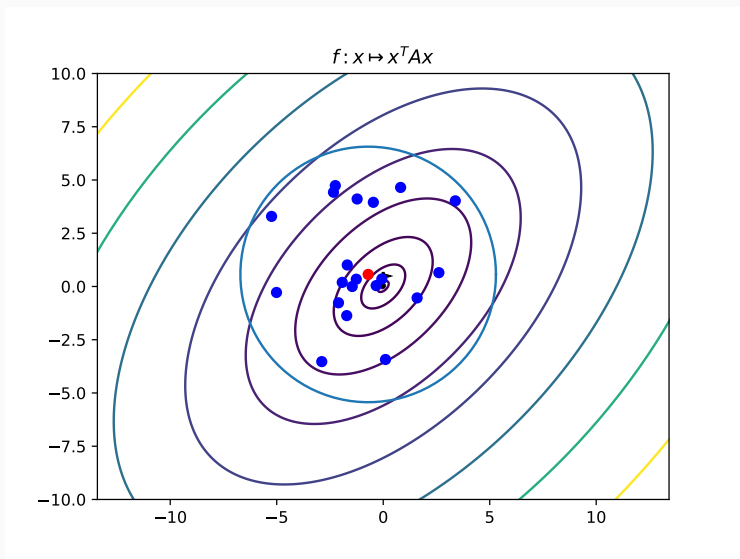
Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ independently



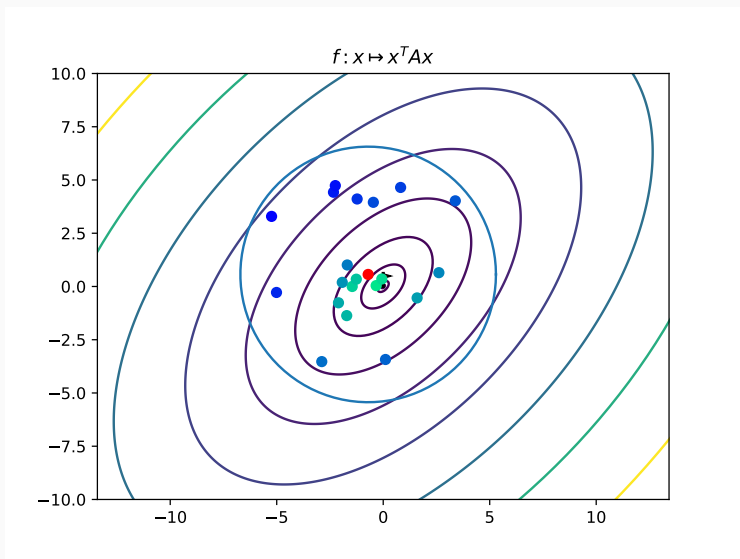
Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ independently



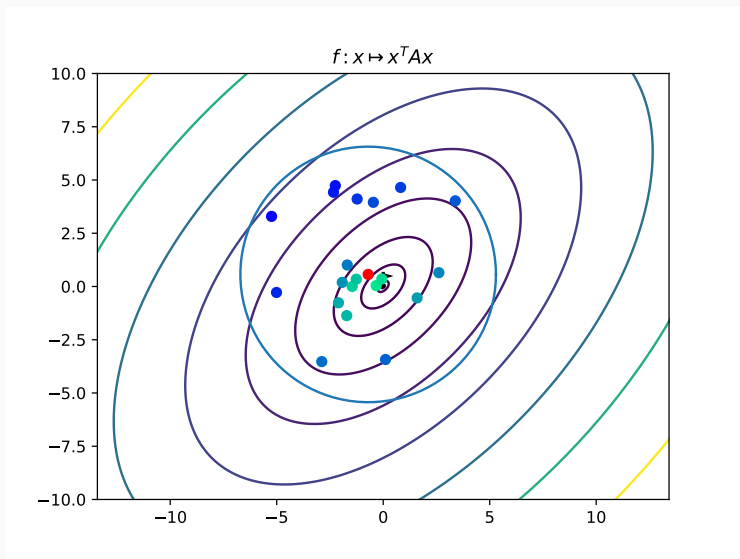
Rank the x_{t+1}^i w.r.t. their f -values $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda})$



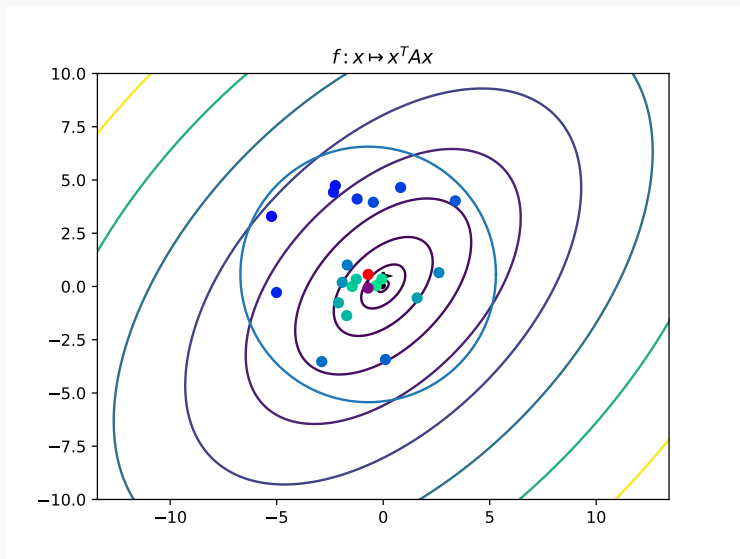
Rank the x_{t+1}^i w.r.t. their f -values $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda})$



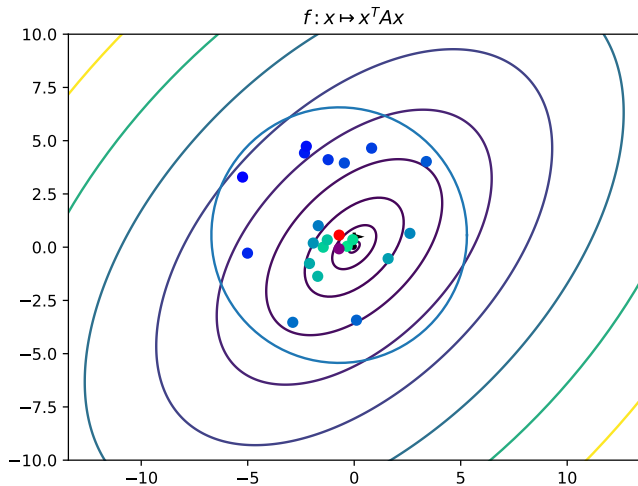
Update the mean $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$



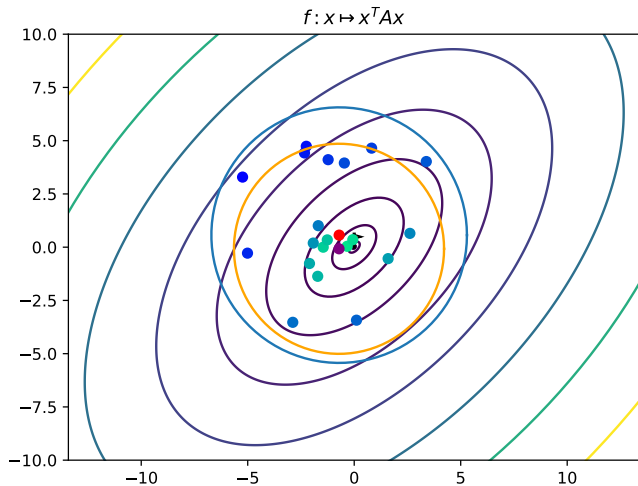
Update the mean $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$



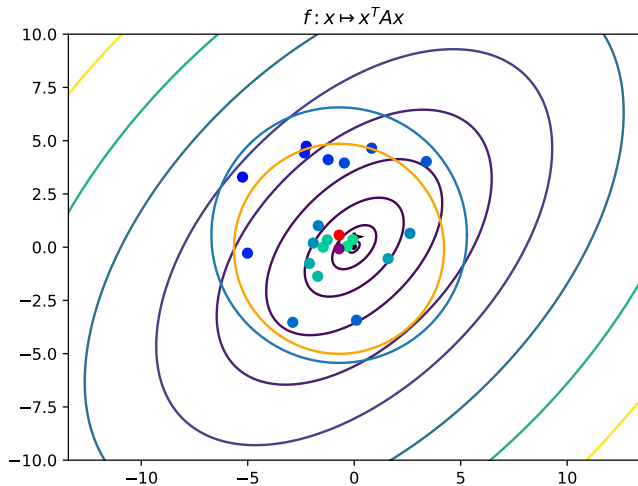
Adapt the stepsize



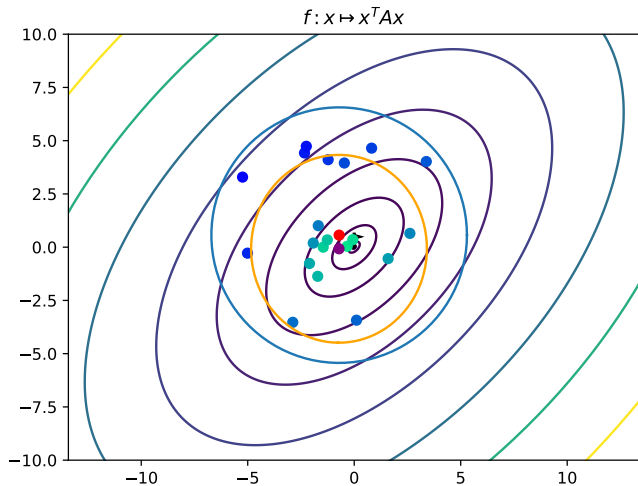
Adapt the stepsize



Adapt the covariance matrix

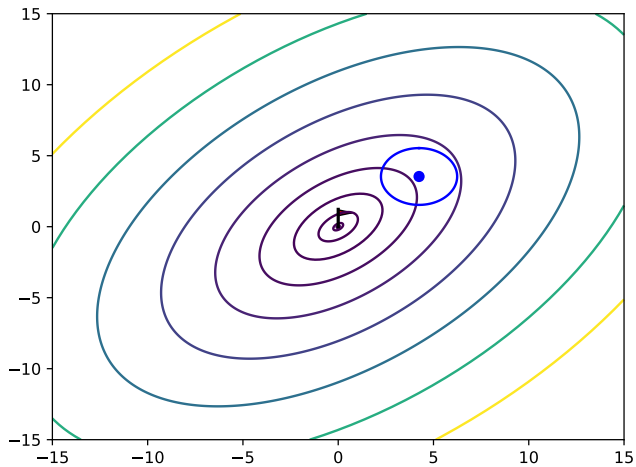


Adapt the covariance matrix

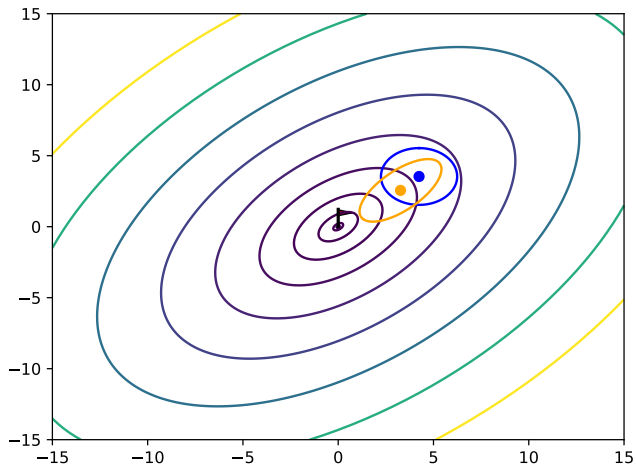


Linear convergence

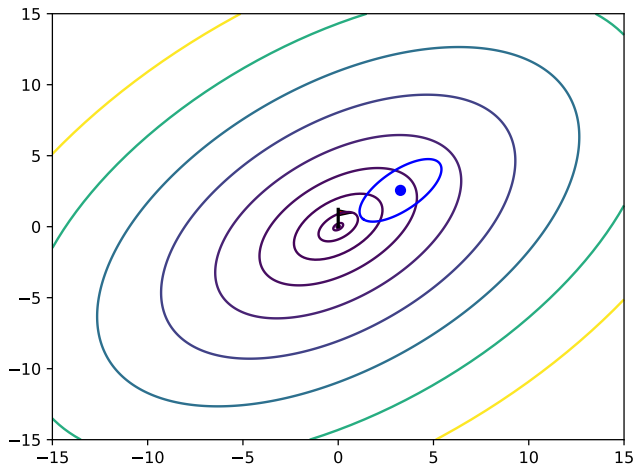
Convergence



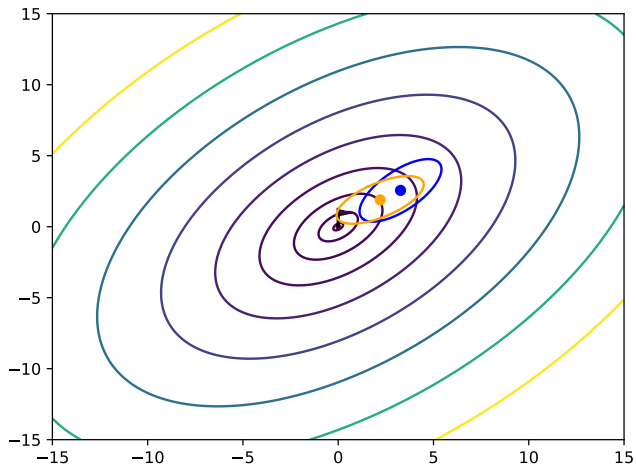
Convergence



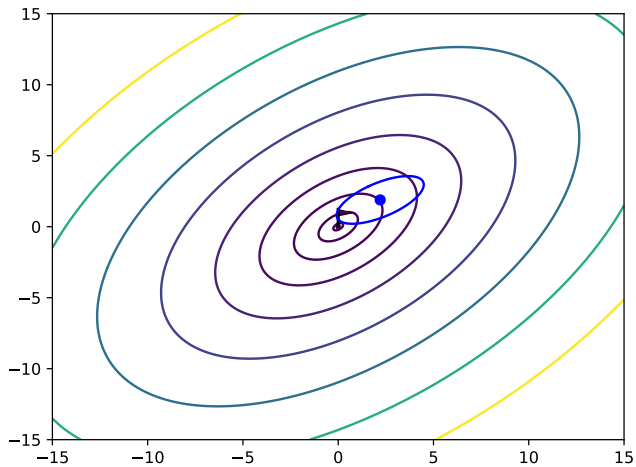
Convergence



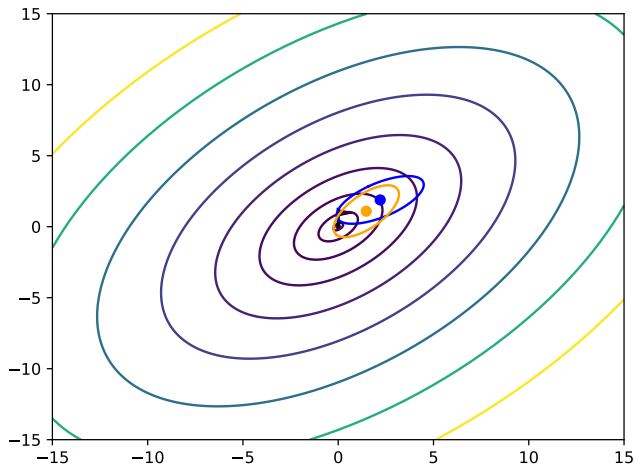
Convergence



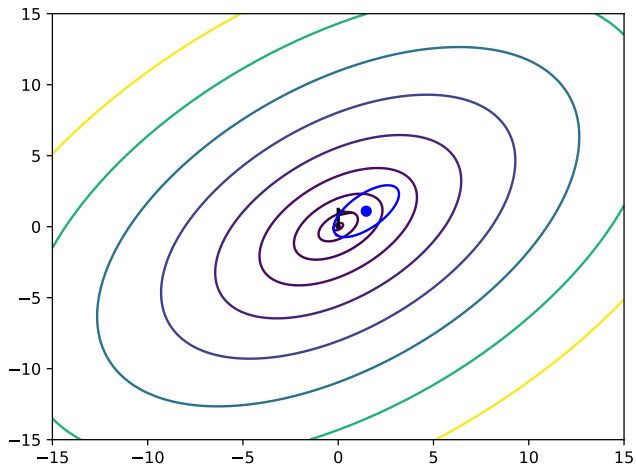
Convergence



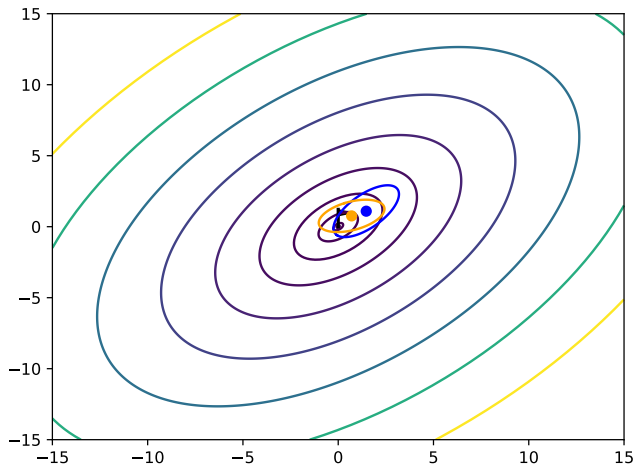
Convergence



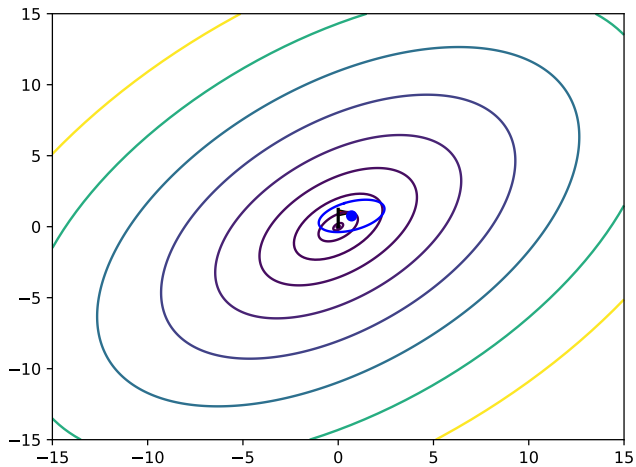
Convergence



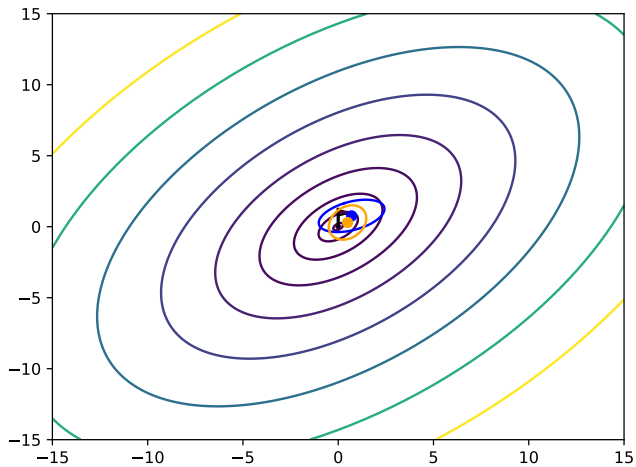
Convergence



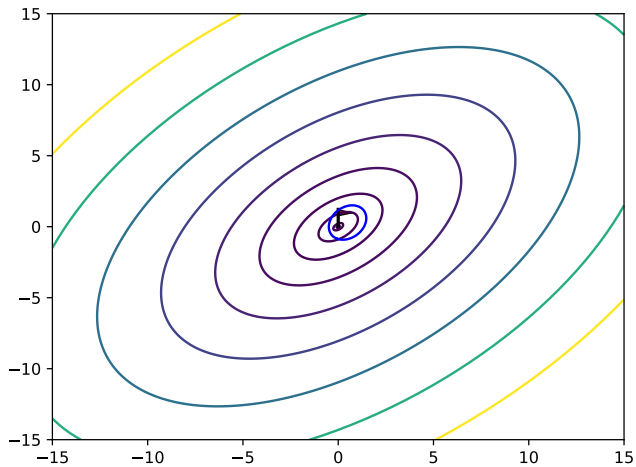
Convergence



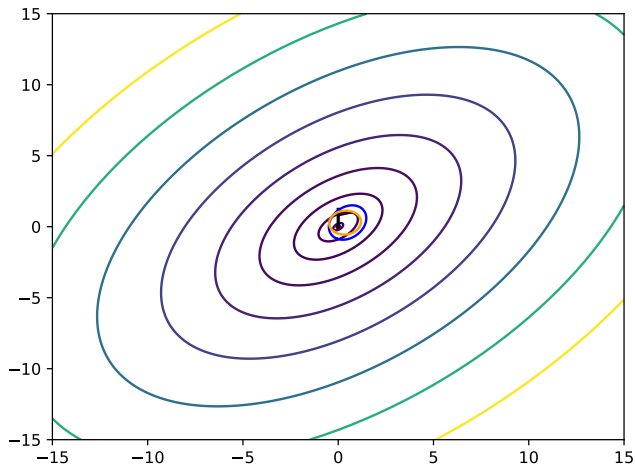
Convergence



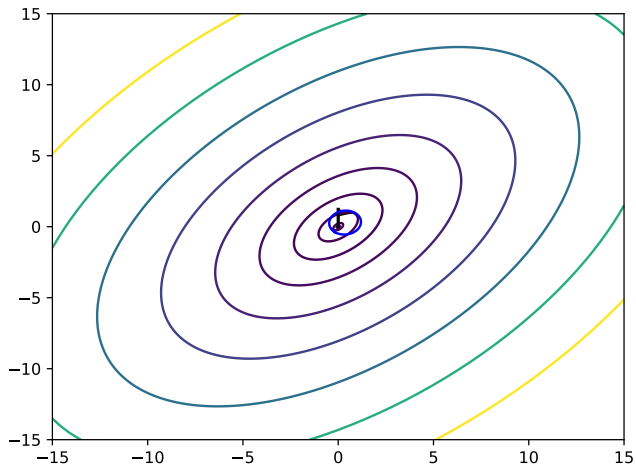
Convergence



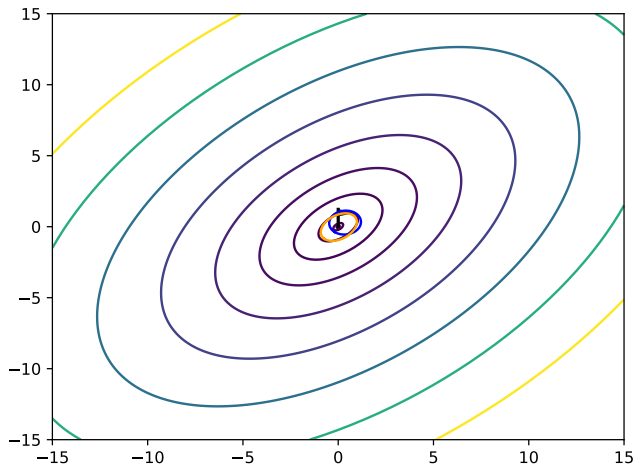
Convergence



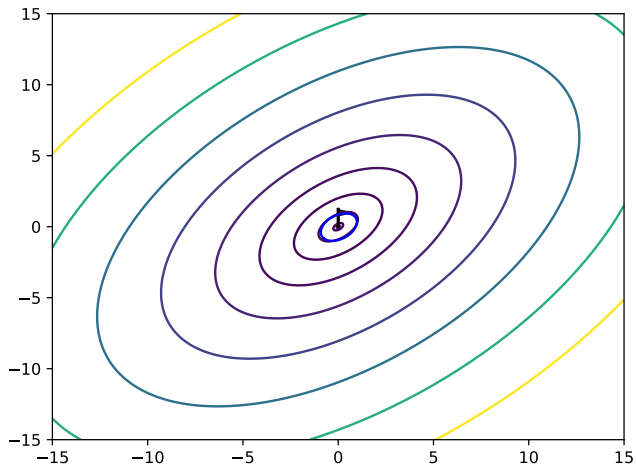
Convergence



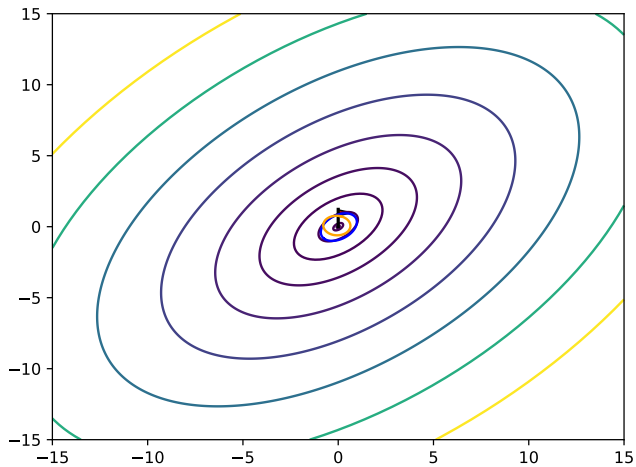
Convergence



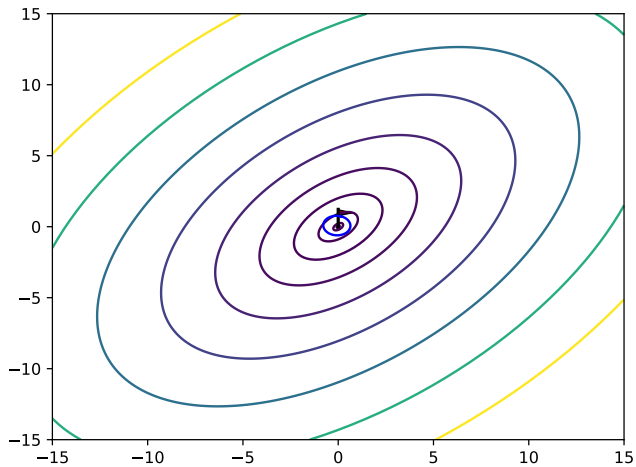
Convergence



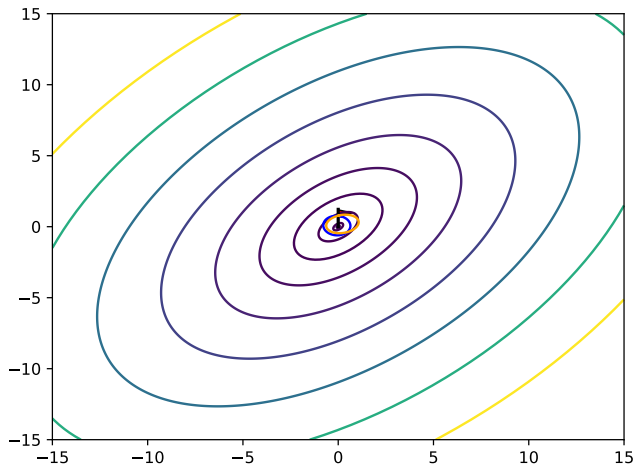
Convergence



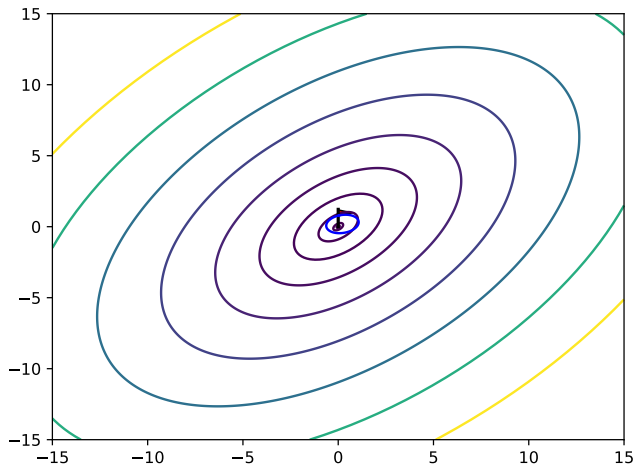
Convergence



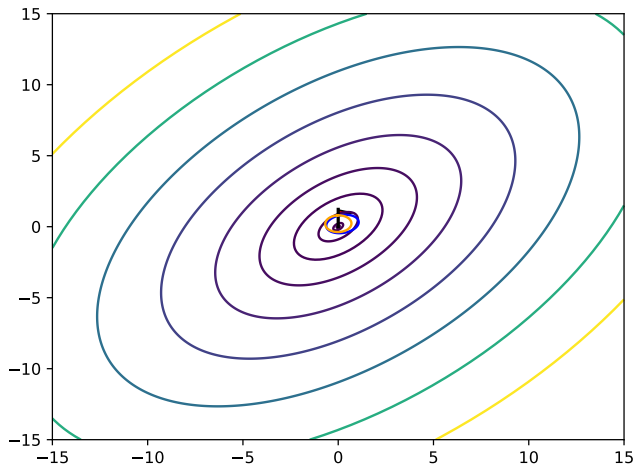
Convergence



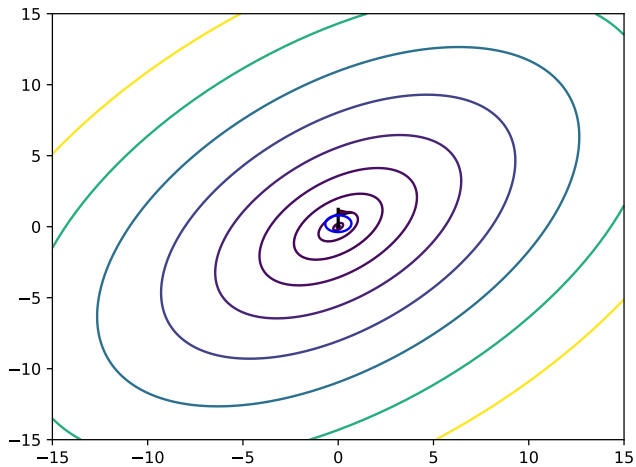
Convergence



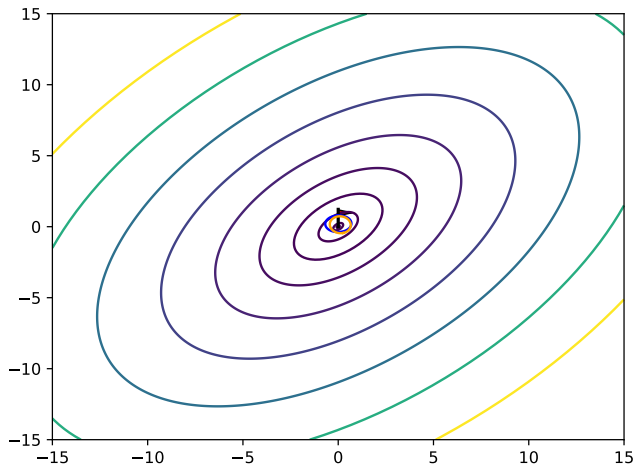
Convergence



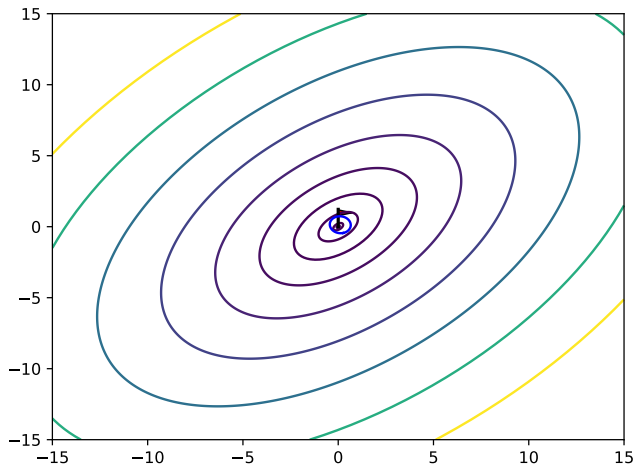
Convergence



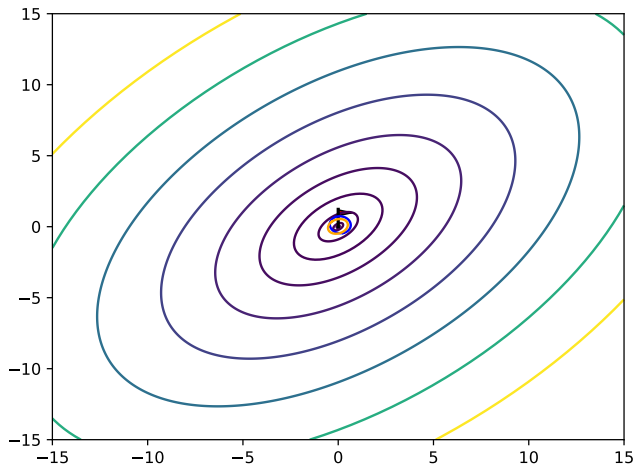
Convergence



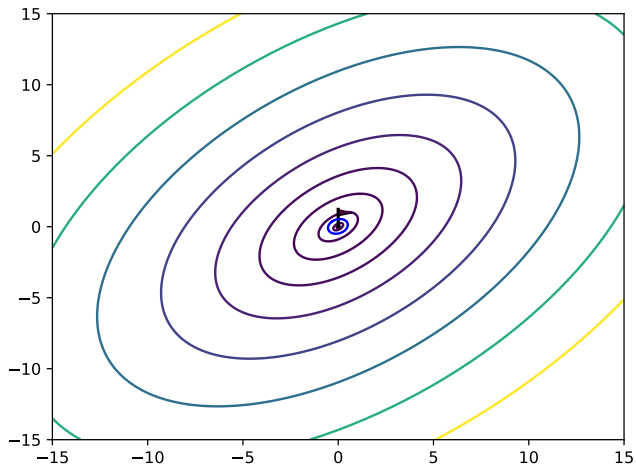
Convergence



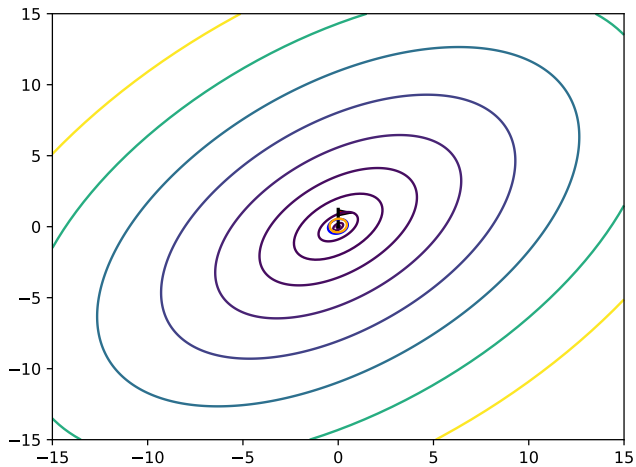
Convergence



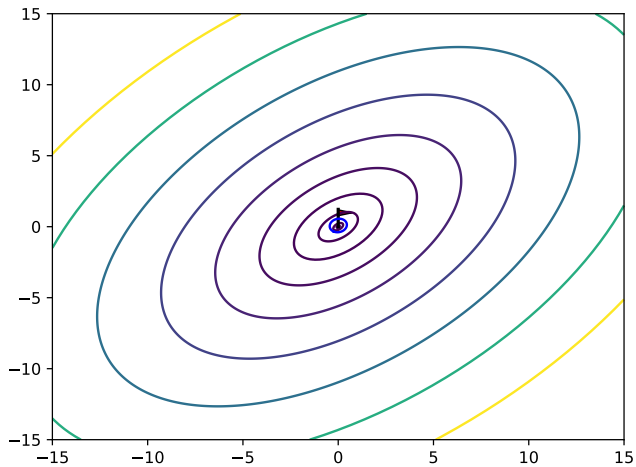
Convergence



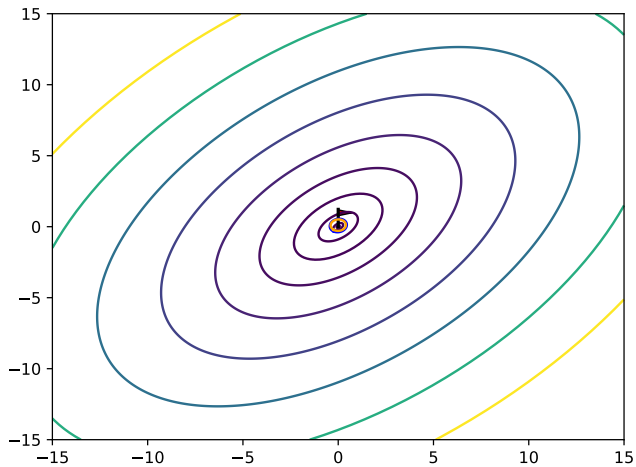
Convergence



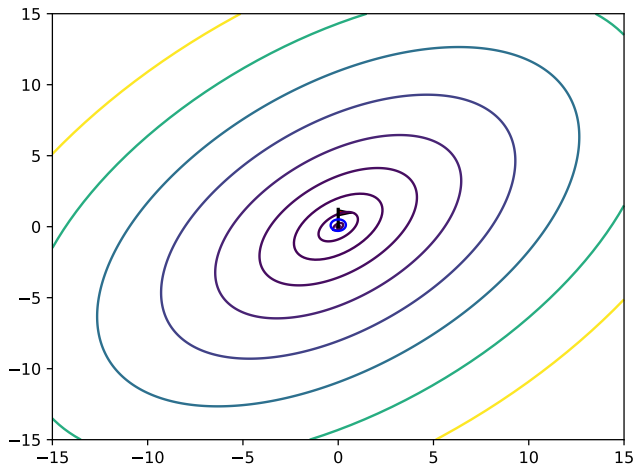
Convergence



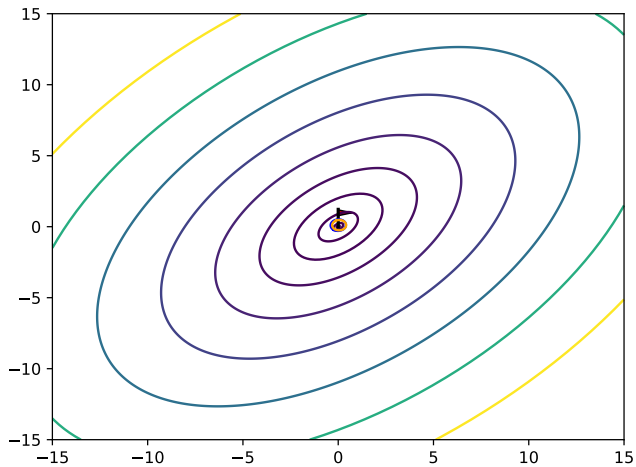
Convergence



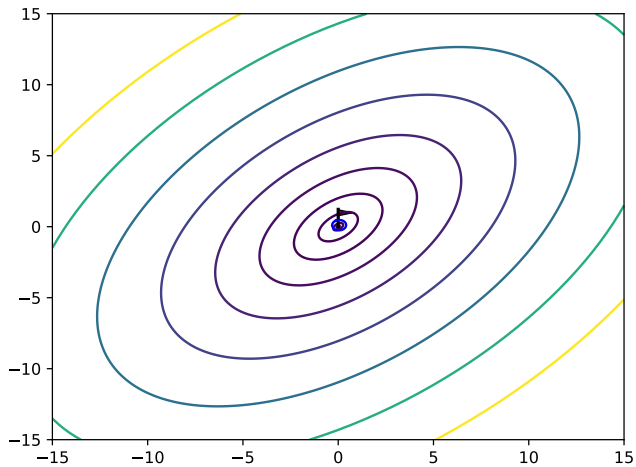
Convergence



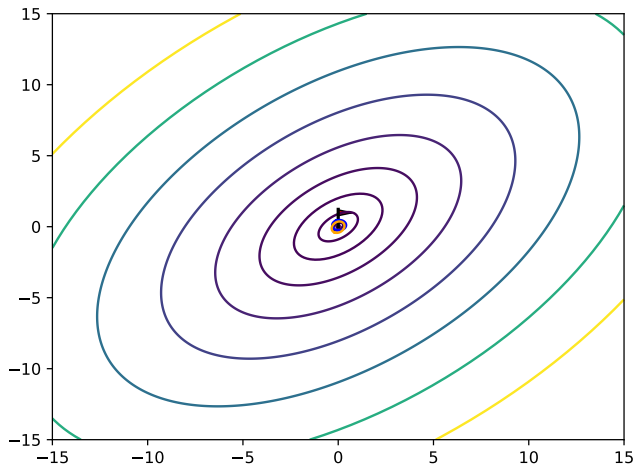
Convergence



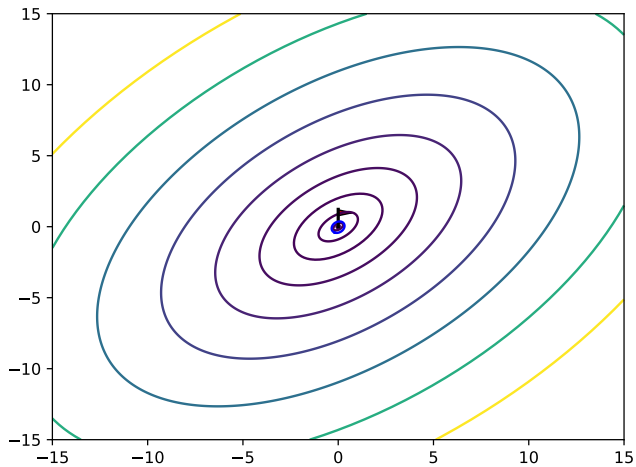
Convergence



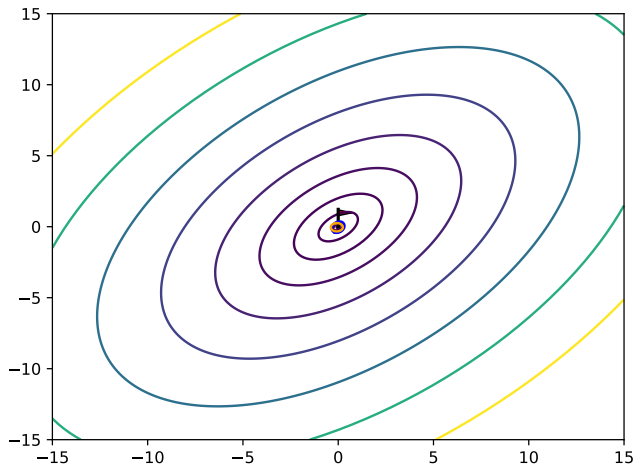
Convergence



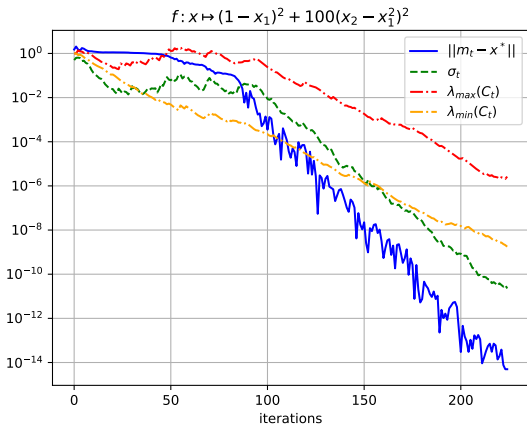
Convergence



Convergence

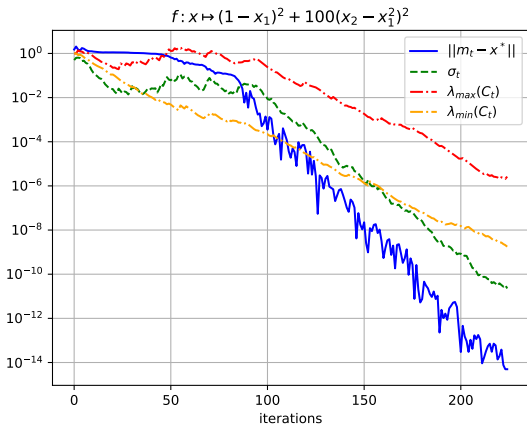


Linear convergence



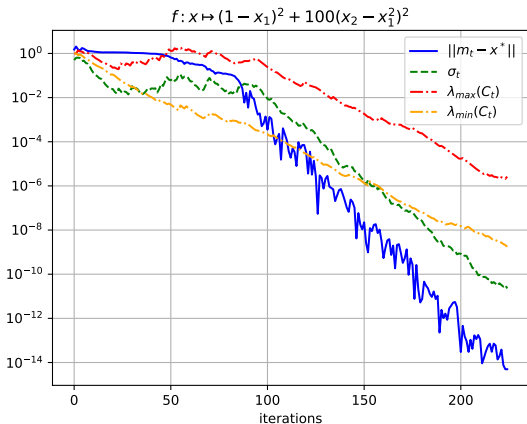
$$\frac{\|m_t - x^*\|}{\|m_0 - x^*\|}$$

Linear convergence



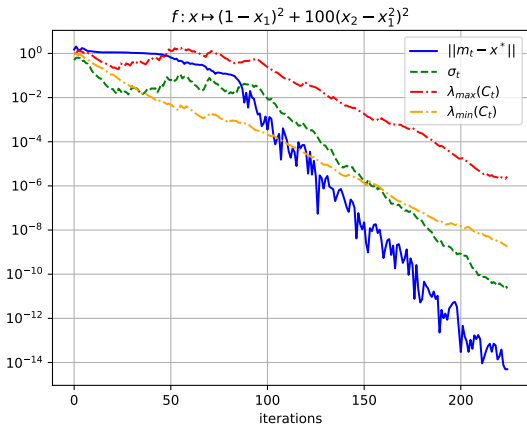
$$\log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|}$$

Linear convergence



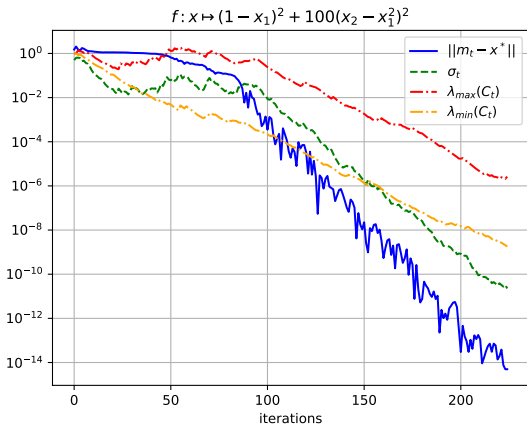
$$\log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = -\text{CR} \times t$$

Linear convergence



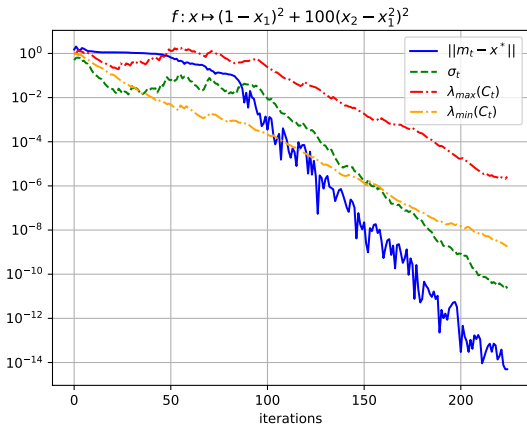
$$\frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = -\text{CR}$$

Linear convergence



$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = -\text{CR}$$

Linear convergence



$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = -CR$$

Summary

- We have **invariance by increasing transformation**

Summary

- We have **invariance by increasing transformation**, i.e. if $g: \mathbb{R} \rightarrow \mathbb{R}$ is **increasing**

Summary

- We have **invariance by increasing transformation**, i.e. if $g: \mathbb{R} \rightarrow \mathbb{R}$ is **increasing**, then **minimizing f and $g \circ f$ will be equivalent**

CMA-ES : presentation

At each iteration $t \in \mathbb{N}$, given a **mean** m_t , a **stepsize** σ_t and a **covariance matrix** C_t :

1. Generate λ offspring $x_{t+1}^i \sim \mathcal{N}(m_t, \sigma_t^2 C_t)$ *independently*
2. Rank the x_{t+1}^i w.r.t. their f -values: $f(x_{t+1}^{1:\lambda}) \leq \dots \leq f(x_{t+1}^{\lambda:\lambda})$
3. Update the mean: $m_{t+1} = \sum w_i x_{t+1}^{i:\lambda}$

The best offspring have the largest weights: $w_1 \geq w_2 \dots$

4. Update the stepsize:
Increase the stepsize if the path taken by the mean is larger than expected (assuming no selection)
5. Update the covariance matrix

$$C_{t+1} = (1 - c)C_t + c \sum w_i \frac{(x_{t+1}^{i:\lambda} - m_t)(x_{t+1}^{i:\lambda} - m_t)^T}{\sigma_t^2}$$

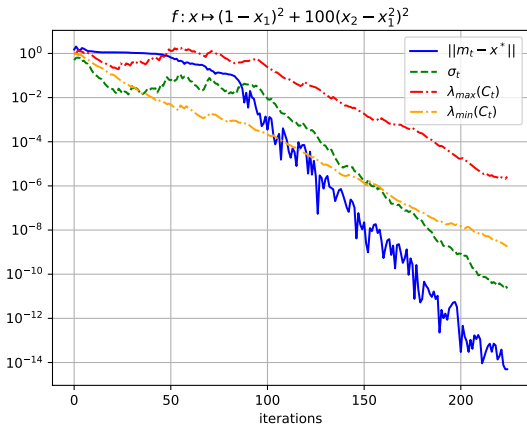
Summary

- We have **invariance by increasing transformation**, i.e. if $g: \mathbb{R} \rightarrow \mathbb{R}$ is **increasing**, then **minimizing f and $g \circ f$ will be equivalent**

Summary

- We have invariance by increasing transformation, i.e. if $g: \mathbb{R} \rightarrow \mathbb{R}$ is increasing, then minimizing f and $g \circ f$ will be equivalent
- We observe **linear convergence** $m_t \rightarrow x^*$

Linear convergence



$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = -CR$$

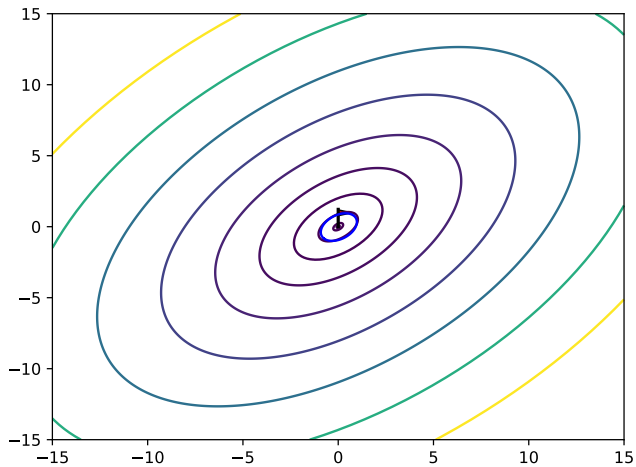
Summary

- We have invariance by increasing transformation, i.e. if $g: \mathbb{R} \rightarrow \mathbb{R}$ is increasing, then minimizing f and $g \circ f$ will be equivalent
- We observe **linear convergence** $m_t \rightarrow x^*$

Summary

- We have invariance by increasing transformation, i.e. if $g: \mathbb{R} \rightarrow \mathbb{R}$ is increasing, then minimizing f and $g \circ f$ will be equivalent
- We observe linear convergence $m_t \rightarrow x^*$
- The covariance matrix **learns the inverse Hessian of f** (when f is e.g. convex-quadratic)

Convergence



Summary

- We have invariance by increasing transformation, i.e. if $g: \mathbb{R} \rightarrow \mathbb{R}$ is increasing, then minimizing f and $g \circ f$ will be equivalent
- We observe linear convergence $m_t \rightarrow x^*$
- The covariance matrix **learns the inverse Hessian of f** (when f is e.g. convex-quadratic)

Summary

- We have invariance by increasing transformation, i.e. if $g: \mathbb{R} \rightarrow \mathbb{R}$ is increasing, then minimizing f and $g \circ f$ will be equivalent
- We observe linear convergence $m_t \rightarrow x^*$
- The covariance matrix learns the inverse Hessian of f (when f is e.g. convex-quadratic)

Goal:

Summary

- We have invariance by increasing transformation, i.e. if $g: \mathbb{R} \rightarrow \mathbb{R}$ is increasing, then minimizing f and $g \circ f$ will be equivalent
- We observe linear convergence $m_t \rightarrow x^*$
- The covariance matrix learns the inverse Hessian of f (when f is e.g. convex-quadratic)

Goal: proof of linear convergence

Summary

- We have invariance by increasing transformation, i.e. if $g: \mathbb{R} \rightarrow \mathbb{R}$ is increasing, then minimizing f and $g \circ f$ will be equivalent
- We observe linear convergence $m_t \rightarrow x^*$
- The covariance matrix learns the inverse Hessian of f (when f is e.g. convex-quadratic)

Goal: proof of linear convergence and of the learning of the inverse Hessian

Analysis via Markov chains

Markov chains

A **Markov chain** is a random sequence $(\phi_t)_{t \in \mathbb{N}}$ such that

$$\text{Distribution}(\phi_{t+1} \mid \phi_0, \dots, \phi_t) = \text{Distribution}(\phi_{t+1} \mid \phi_t)$$

A **Markov chain** is a random sequence $(\phi_t)_{t \in \mathbb{N}}$ such that

$$\text{Distribution}(\phi_{t+1} \mid \phi_0, \dots, \phi_t) = \text{Distribution}(\phi_{t+1} \mid \phi_t)$$

- The Markov chain $(\phi_t)_{t \in \mathbb{N}}$ is **irreducible** if all states¹ are reachable in finite time with positive probability.

¹with positive measure (e.g. Lebesgue measure)

A **Markov chain** is a random sequence $(\phi_t)_{t \in \mathbb{N}}$ such that

$$\text{Distribution}(\phi_{t+1} \mid \phi_0, \dots, \phi_t) = \text{Distribution}(\phi_{t+1} \mid \phi_t)$$

- The Markov chain $(\phi_t)_{t \in \mathbb{N}}$ is **irreducible** if all states¹ are reachable in finite time with positive probability.
- A probability distribution π is **invariant** for the Markov chain $(\phi_t)_{t \in \mathbb{N}}$ when

$$\phi_t \sim \pi \Rightarrow \phi_{t+1} \sim \pi$$

¹with positive measure (e.g. Lebesgue measure)

A **Markov chain** is a random sequence $(\phi_t)_{t \in \mathbb{N}}$ such that

$$\text{Distribution}(\phi_{t+1} \mid \phi_0, \dots, \phi_t) = \text{Distribution}(\phi_{t+1} \mid \phi_t)$$

- The Markov chain $(\phi_t)_{t \in \mathbb{N}}$ is **irreducible** if all states¹ are reachable in finite time with positive probability.
- A probability distribution π is **invariant** for the Markov chain $(\phi_t)_{t \in \mathbb{N}}$ when

$$\phi_t \sim \pi \Rightarrow \phi_{t+1} \sim \pi$$

(π is a "fixed point" for $(\phi_t)_{t \in \mathbb{N}}$)

¹with positive measure (e.g. Lebesgue measure)

A **Markov chain** is a random sequence $(\phi_t)_{t \in \mathbb{N}}$ such that

$$\text{Distribution}(\phi_{t+1} \mid \phi_0, \dots, \phi_t) = \text{Distribution}(\phi_{t+1} \mid \phi_t)$$

- The Markov chain $(\phi_t)_{t \in \mathbb{N}}$ is **irreducible** if all states¹ are reachable in finite time with positive probability.
- A probability distribution π is **invariant** for the Markov chain $(\phi_t)_{t \in \mathbb{N}}$ when

$$\phi_t \sim \pi \Rightarrow \phi_{t+1} \sim \pi$$

(π is a "fixed point" for $(\phi_t)_{t \in \mathbb{N}}$)

- The Markov chain is **ergodic** when it satisfies the following LLN

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} f(\phi_t) = \int f(x) \pi(dx).$$

¹with positive measure (e.g. Lebesgue measure)

CMA-ES as a Markov chain

Consider the random sequence

$$\left(\underbrace{m_t}_{\text{mean}}, \underbrace{\sigma_t}_{\text{stepsize}}, \underbrace{C_t}_{\text{covariance matrix}} \right)$$

CMA-ES as a Markov chain

Consider the random sequence

$$\left(\underbrace{m_t}_{\text{mean}}, \underbrace{\sigma_t}_{\text{stepsize}}, \underbrace{C_t}_{\text{covariance matrix}} \right)$$

This defines a Markov chain!

CMA-ES as a Markov chain

Consider the random sequence

$$\left(\underbrace{m_t}_{\text{mean}}, \underbrace{\sigma_t}_{\text{stepsize}}, \underbrace{C_t}_{\text{covariance matrix}} \right)$$

This defines a Markov chain!

Question: **Could we use the LLN for Markov chains to prove linear convergence for CMA-ES?**

Invariant measure for CMA-ES?

Suppose that $(m_t, \sigma_t, C_t)_{t \in \mathbb{N}}$ has an invariant measure π .

Invariant measure for CMA-ES?

Suppose that $(m_t, \sigma_t, C_t)_{t \in \mathbb{N}}$ has an invariant measure π .

$$(m_t, \sigma_t, C_t) \sim \pi \Rightarrow (m_{t+1}, \sigma_{t+1}, C_{t+1}) \sim \pi$$

Invariant measure for CMA-ES?

Suppose that $(m_t, \sigma_t, C_t)_{t \in \mathbb{N}}$ has an invariant measure π .

$$(m_t, \sigma_t, C_t) \sim \pi \Rightarrow (m_{t+1}, \sigma_{t+1}, C_{t+1}) \sim \pi$$

We do not progress towards the optimum anymore!

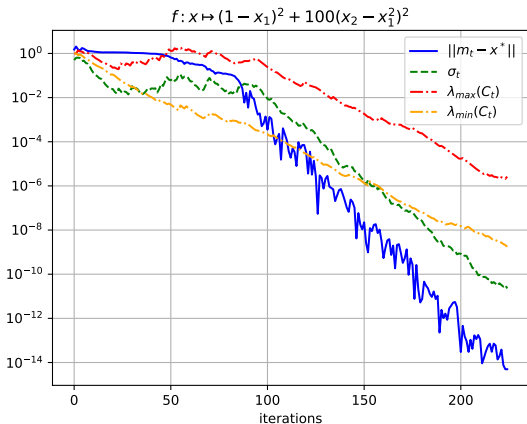
Invariant measure for CMA-ES?

Suppose that $(m_t, \sigma_t, C_t)_{t \in \mathbb{N}}$ has an invariant measure π .

$$(m_t, \sigma_t, C_t) \sim \pi \Rightarrow (m_{t+1}, \sigma_{t+1}, C_{t+1}) \sim \pi$$

We do not progress towards the optimum anymore! The **existence of an invariant measure** seems to be **incompatible** with the **convergence** to the optimum.

Linear convergence



$$\lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{\|m_t - x^*\|}{\|m_0 - x^*\|} = -CR$$

$$\|m_t - x^*\|, \sigma_t \text{ and } \lambda_{\min}(C_t) \rightarrow 0$$

Normalization

$\|m_t - x^*\|$, σ_t and $\lambda_{\min}(C_t) \rightarrow 0$

$$Z_t \stackrel{\text{def}}{=} \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

Normalization

$\|m_t - x^*\|$, σ_t and $\lambda_{\min}(C_t) \rightarrow 0$

$$Z_t \stackrel{\text{def}}{=} \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

The sequence $(Z_t)_{t \in \mathbb{N}}$ could eventually become **stationary**

$\|m_t - x^*\|$, σ_t and $\lambda_{\min}(C_t) \rightarrow 0$

$$Z_t \stackrel{\text{def}}{=} \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

The sequence $(Z_t)_{t \in \mathbb{N}}$ could eventually become stationary

Proposition (Normalized Markov chain)

The sequence

$$\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)_{t \in \mathbb{N}}$$

defines a Markov chain.

Normalization

$\|m_t - x^*\|$, σ_t and $\lambda_{\min}(C_t) \rightarrow 0$

$$Z_t \stackrel{\text{def}}{=} \frac{m_t - x^*}{\sigma_t \sqrt{\lambda_{\min}(C_t)}}$$

The sequence $(Z_t)_{t \in \mathbb{N}}$ could eventually become stationary

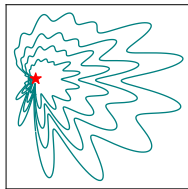
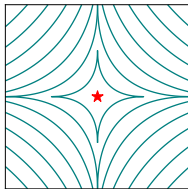
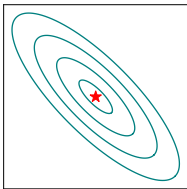
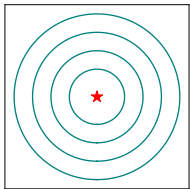
Proposition (Normalized Markov chain)

The sequence

$$\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)_{t \in \mathbb{N}}$$

defines a Markov chain. (if f is **scaling-invariant**)

Scaling-invariant functions

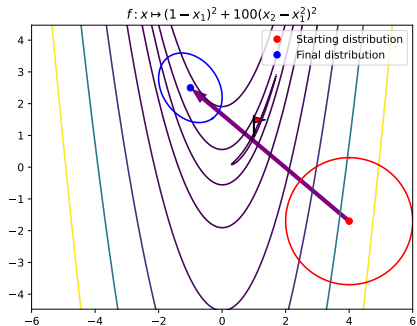


Irreducibility of CMA-ES

A Markov chain $(\phi_t)_{t \in \mathbb{N}}$ is **irreducible** if all states are reachable in finite time with positive probability.

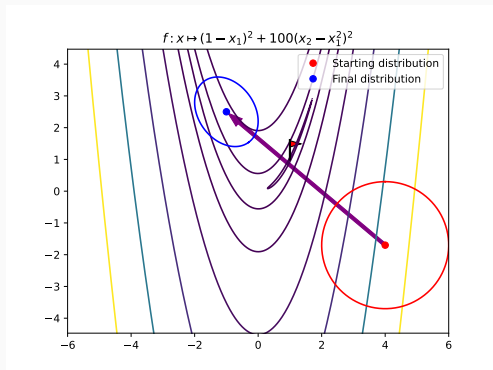
Irreducibility of CMA-ES

A Markov chain $(\phi_t)_{t \in \mathbb{N}}$ is **irreducible** if all states are reachable in finite time with positive probability.



Irreducibility of CMA-ES

A Markov chain $(\phi_t)_{t \in \mathbb{N}}$ is **irreducible** if all states are reachable in finite time with positive probability.



Given a starting and a final distributions, **can we reach the final distribution in finite time with positive probability?**

Theorem (Irreducibility of the normalized chain)

When minimizing a scaling-invariant function f with Lebesgue-negligible level sets, the sequence

$$\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)_{t \in \mathbb{N}}$$

defines a irreducible, aperiodic Markov chain, and compact sets are small sets.

Ergodicity of the normalized chain²

$$\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)_{t \in \mathbb{N}}$$

is **ergodic** if

²Sean P. Meyn and Richard L. Tweedie. **Markov Chains and Stochastic Stability**. Springer Science & Business Media, Dec. 2012. ISBN: 978-1-4471-3267-7.

Ergodicity of the normalized chain²

$$\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)_{t \in \mathbb{N}}$$

is **ergodic** (and satisfies a LLN) if

²Sean P. Meyn and Richard L. Tweedie. **Markov Chains and Stochastic Stability**. Springer Science & Business Media, Dec. 2012. ISBN: 978-1-4471-3267-7.

Ergodicity of the normalized chain²

$$\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)_{t \in \mathbb{N}}$$

is **ergodic** (and satisfies a LLN) if

- it is **irreducible** and aperiodic

²Sean P. Meyn and Richard L. Tweedie. **Markov Chains and Stochastic Stability**. Springer Science & Business Media, Dec. 2012. ISBN: 978-1-4471-3267-7.

Ergodicity of the normalized chain²

$$\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)_{t \in \mathbb{N}}$$

is **ergodic** (and satisfies a LLN) if

- it is irreducible and aperiodic
- it satisfies the following **drift** condition: $\exists V: X \rightarrow [0, +\infty]$

$$\mathbb{E}_t \left[V \left(Z_{t+1}, \frac{C_{t+1}}{\lambda_{\min}(C_{t+1})} \right) \right] \leq (1 - \varepsilon) V \left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)$$

²Sean P. Meyn and Richard L. Tweedie. **Markov Chains and Stochastic Stability**. Springer Science & Business Media, Dec. 2012. ISBN: 978-1-4471-3267-7.

Ergodicity of the normalized chain²

$$\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)_{t \in \mathbb{N}}$$

is **ergodic** (and satisfies a LLN) if

- it is irreducible and aperiodic
- it satisfies the following drift condition: $\exists V: X \rightarrow [0, +\infty]$

$$\mathbb{E}_t \left[V \left(Z_{t+1}, \frac{C_{t+1}}{\lambda_{\min}(C_{t+1})} \right) \right] \leq (1 - \varepsilon) V \left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)$$

for $\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)$ **outside of a compact** \mathcal{K} .

²Sean P. Meyn and Richard L. Tweedie. **Markov Chains and Stochastic Stability**. Springer Science & Business Media, Dec. 2012. ISBN: 978-1-4471-3267-7.

Ergodicity of the normalized chain²

$$\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)_{t \in \mathbb{N}}$$

is **ergodic** (and satisfies a LLN) if

- it is irreducible and aperiodic
- it satisfies the following drift condition: $\exists V: X \rightarrow [0, +\infty]$

$$\mathbb{E}_t \left[V \left(Z_{t+1}, \frac{C_{t+1}}{\lambda_{\min}(C_{t+1})} \right) \right] \leq (1 - \varepsilon) V \left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)$$

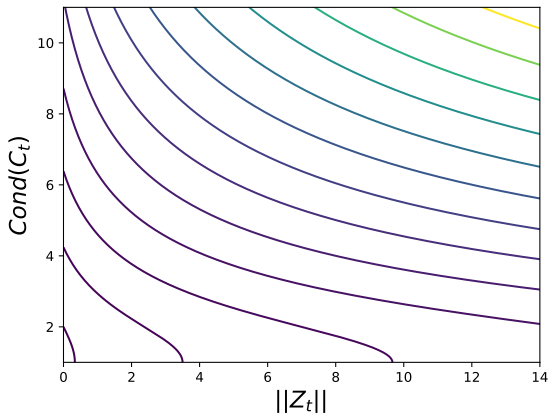
for $\left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)$ outside of a compact \mathcal{K} .

The function V is called the **potential function** or the **drift function** or the **Lyapunov function**.

²Sean P. Meyn and Richard L. Tweedie. **Markov Chains and Stochastic Stability**. Springer Science & Business Media, Dec. 2012.

Drift condition

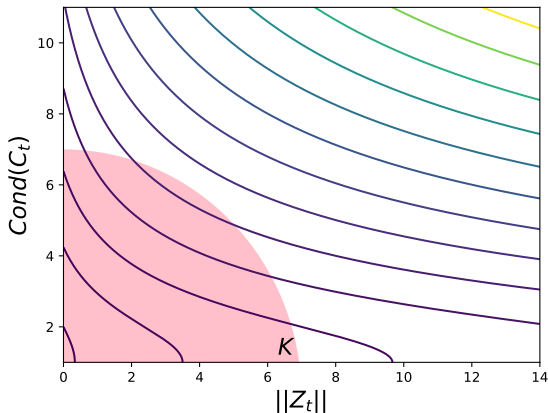
$$\mathbb{E}_t \left[V \left(Z_{t+1}, \frac{C_{t+1}}{\lambda_{\min}(C_{t+1})} \right) \right] \leq (1 - \varepsilon) V \left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)$$



Drift condition

$$\mathbb{E}_t \left[V \left(Z_{t+1}, \frac{C_{t+1}}{\lambda_{\min}(C_{t+1})} \right) \right] \leq (1 - \varepsilon) V \left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)$$

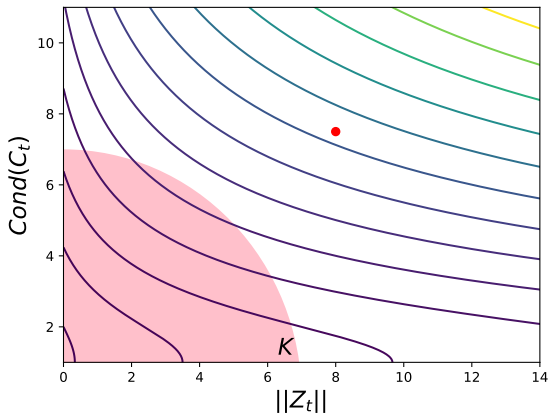
outside of a compact K



Drift condition

$$\mathbb{E}_t \left[V \left(Z_{t+1}, \frac{C_{t+1}}{\lambda_{\min}(C_{t+1})} \right) \right] \leq (1 - \varepsilon) V \left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)$$

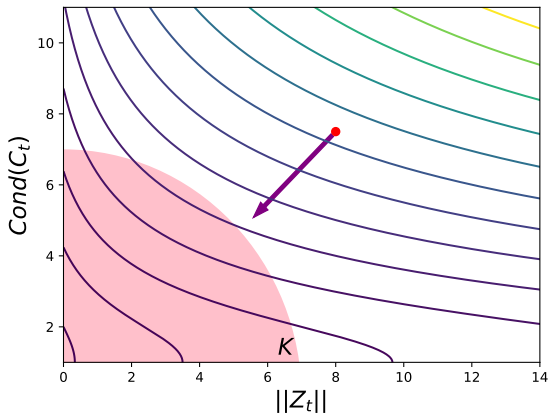
outside of a compact K



Drift condition

$$\mathbb{E}_t \left[V \left(Z_{t+1}, \frac{C_{t+1}}{\lambda_{\min}(C_{t+1})} \right) \right] \leq (1 - \varepsilon) V \left(Z_t, \frac{C_t}{\lambda_{\min}(C_t)} \right)$$

outside of a compact K



Theorem (Drift condition for the normalized chain)

When minimizing a spherical function $f : x \mapsto g(x^T x)$ ($g : \mathbb{R} \rightarrow \mathbb{R}$ increasing), then the irreducible, aperiodic Markov chain $(Z_t, C_t / \lambda_{\min}(C_t))_{t \in \mathbb{N}}$ satisfies a Foster-Lyapunov condition with the potential defined by

$$V(Z, C) = \sum_{k=1}^d \left\{ \frac{\lambda_k(C)}{\lambda_1(C)} |\langle v_k(C), Z \rangle|^2 \right\} + \beta \times \text{Cond}(C)$$

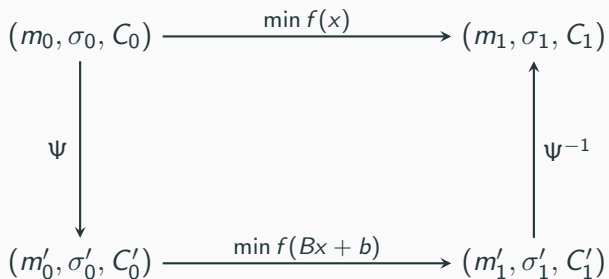
Theorem (Drift condition for the normalized chain)

When minimizing a spherical function $f : x \mapsto g(x^T x)$ ($g : \mathbb{R} \rightarrow \mathbb{R}$ increasing), then the irreducible, aperiodic Markov chain $(Z_t, C_t / \lambda_{\min}(C_t))_{t \in \mathbb{N}}$ satisfies a Foster-Lyapunov condition with the potential defined by

$$V(Z, C) = \sum_{k=1}^d \left\{ \frac{\lambda_k(C)}{\lambda_1(C)} |\langle v_k(C), Z \rangle|^2 \right\} + \beta \times \text{Cond}(C)$$

This can be generalized to when minimizing ellipsoid functions $f(x) = g(x^T Hx)$ using the **affine-invariance** of CMA-ES.

Affine-Invariance



Proof of linear convergence

$$\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} \rightarrow -\text{CR?}$$

Proof of linear convergence

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

Proof of linear convergence

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

Proof of linear convergence

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

Proof of linear convergence

$$\begin{aligned}\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\sigma_{t+1} \sqrt{\lambda_{\min}(C_{t+1})}} + \log \sigma_{t+1} \\ &\quad + \frac{1}{2} \log \lambda_{\min}(C_{t+1}) - \log \|m_t - x^*\|\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \log \|Z_{t+1}\| + \log \sigma_{t+1} \\ &\quad + \frac{1}{2} \log \lambda_{\min}(C_{t+1}) - \log \|m_t - x^*\|\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \log \|Z_{t+1}\| + \log \sigma_{t+1} \\ &\quad + \frac{1}{2} \log \lambda_{\min}(C_{t+1}) - \log \|m_t - x^*\|\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \log \|Z_{t+1}\| + \log \sigma_{t+1} \\ &\quad + \frac{1}{2} \log \lambda_{\min}(C_{t+1}) - \log \|m_t - x^*\|\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \log \|Z_{t+1}\| - \log \|Z_t\| \\ &\quad + \log \frac{\sigma_{t+1}}{\sigma_t} + \frac{1}{2} \log \frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \log \|Z_{t+1}\| - \log \|Z_t\| \\ &\quad + \log \frac{\sigma_{t+1}}{\sigma_t} + \frac{1}{2} \log \frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \|Z_{t+1}\| - \log \|Z_t\| \\ &\quad + \log \frac{\sigma_{t+1}}{\sigma_t} + \frac{1}{2} \log \frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \|Z_{t+1}\| - \log \|Z_t\| \\ &\quad + \log \frac{\sigma_{t+1}}{\sigma_t} + \frac{1}{2} \log \frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \int \log \|z\| d\pi - \int \log \|z\| d\pi \\ &\quad + \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\sigma_{t+1}}{\sigma_t} + \frac{1}{2} \log \frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\sigma_{t+1}}{\sigma_t} + \frac{1}{2} \log \frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\sigma_{t+1}}{\sigma_t} + \frac{1}{2} \log \frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\sigma_T}{\sigma_0} + \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{2} \log \frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\sigma_T}{\sigma_0} + \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{2} \log \frac{\lambda_{\min}(C_{t+1})}{\lambda_{\min}(C_t)}\end{aligned}$$

Proof of linear convergence

$$\begin{aligned}\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} &= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \log \frac{\|m_{t+1} - x^*\|}{\|m_t - x^*\|} \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\sigma_T}{\sigma_0} + \frac{1}{2T} \log \frac{\lambda_{\min}(C_T)}{\lambda_{\min}(C_0)}\end{aligned}$$

Proof of linear convergence

$$\lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\|m_T - x^*\|}{\|m_0 - x^*\|} = \lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\sigma_T}{\sigma_0} + \frac{1}{2} \lim_{T \rightarrow \infty} \frac{1}{T} \log \frac{\lambda_{\min}(C_T)}{\lambda_{\min}(C_0)}$$

- **Proof of irreducibility of a normalized chain**

- Proof of irreducibility of a normalized chain
- **Found a potential function on which we have a drift condition for ergodicity**

- Proof of irreducibility of a normalized chain
- Found a potential function on which we have a drift condition for ergodicity³

³currently when the objective function is an increasing transformation of a convex-quadratic function

- Proof of irreducibility of a normalized chain
- Found a potential function on which we have a drift condition for ergodicity³
- **Proof of linear convergence**

³currently when the objective function is an increasing transformation of a convex-quadratic function

- Proof of irreducibility of a normalized chain
- Found a potential function on which we have a drift condition for ergodicity³
- Proof of linear convergence
- **When minimizing a convex-quadratic function**

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

³currently when the objective function is an increasing transformation of a convex-quadratic function

Thank you!