

Natural Evolution Strategies for Direct Search

PGMO-COPI 2014
Recent Advances on Continuous
Randomized black-box optimization

Thursday October 30, 2014

Tobias Glasmachers
Institut für Neuroinformatik
Ruhr-Universität Bochum, Germany



INSTITUT
FÜR
NEUROINFORMATIK

RUHR
UNIVERSITÄT
BOCHUM

RUB

Introduction

“Function value free” black-box optimization:

“Function value free” black-box optimization:

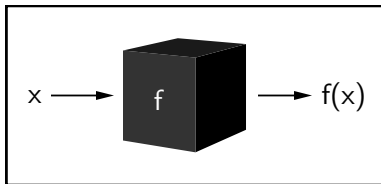
- 1 Minimize $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

“Function value free” black-box optimization:

- 1 Minimize $f : \mathbb{R}^d \rightarrow \mathbb{R}$.
- 2 0-th order (direct search) setting.

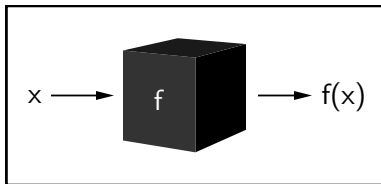
“Function value free” black-box optimization:

- 1 Minimize $f : \mathbb{R}^d \rightarrow \mathbb{R}$.
- 2 0-th order (direct search) setting.
- 3 Black-box model



“Function value free” black-box optimization:

- 1 Minimize $f : \mathbb{R}^d \rightarrow \mathbb{R}$.
- 2 0-th order (direct search) setting.
- 3 Black-box model



- 4 Don't ever rely on function values, only comparisons $f(x) < f(x')$, e.g., for ranking.

Algorithm operations:

Algorithm operations:

- 1 generate solution $x \in \mathbb{R}^d$ – most of the logic

Algorithm operations:

- 1 generate solution $x \in \mathbb{R}^d$ – most of the logic
- 2 evaluate $f(x)$ – black box, most of the computation time

Algorithm operations:

- 1 generate solution $x \in \mathbb{R}^d$ – most of the logic
- 2 evaluate $f(x)$ – black box, most of the computation time
- 3 compare (rank) against other solutions: $f(x) < f(x')$?

Introduction: Challenges

Facing a black-box objective, anything can happen.
We'd better prepare for the following **challenges**:

Facing a black-box objective, anything can happen.
We'd better prepare for the following **challenges**:

- non-smooth or even discontinuous objective

Facing a black-box objective, anything can happen.
We'd better prepare for the following **challenges**:

- non-smooth or even discontinuous objective
- multi-modality

Facing a black-box objective, anything can happen.

We'd better prepare for the following **challenges**:

- non-smooth or even discontinuous objective
- multi-modality
- observations of objective values perturbed by noise

Facing a black-box objective, anything can happen.

We'd better prepare for the following **challenges**:

- non-smooth or even discontinuous objective
- multi-modality
- observations of objective values perturbed by noise
- high dimensionality, e.g., $d \gg 1000$

Facing a black-box objective, anything can happen.

We'd better prepare for the following **challenges**:

- non-smooth or even discontinuous objective
- multi-modality
- observations of objective values perturbed by noise
- high dimensionality, e.g., $d \gg 1000$
- black-box constraints (possibly non-smooth, ...)

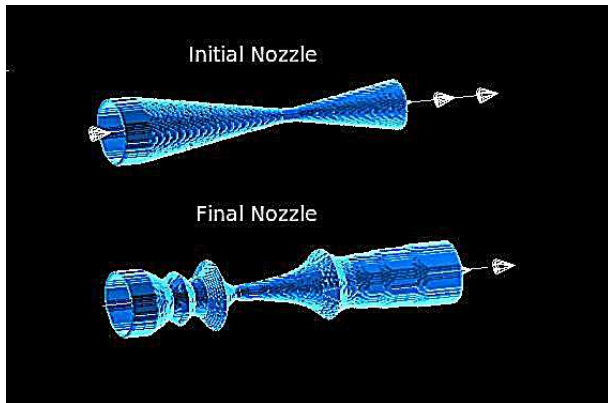
Introduction: Application Examples

Introduction: Application Examples

- 1 Tuning of parameters of expensive simulations, e.g., in engineering.

Introduction: Application Examples

- 1 Tuning of parameters of expensive simulations, e.g., in engineering.

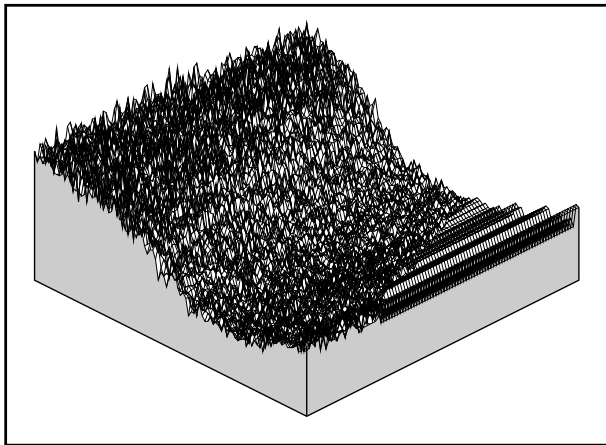


Introduction: Application Examples

- ① Tuning of parameters of expensive simulations, e.g., in engineering.
- ② Tuning of parameters of machine learning algorithms.

Introduction: Application Examples

- 1 Tuning of parameters of expensive simulations, e.g., in engineering.
- 2 Tuning of parameters of machine learning algorithms.



Evolution Strategies

input $m \in \mathbb{R}^d$, $\sigma > 0$, $C(= I) \in \mathbb{R}^{d \times d}$

loop

sample “offspring” $x_1, \dots, x_\lambda \sim \mathcal{N}(m, \sigma^2 C)$

evaluate $f(x_1), \dots, f(x_\lambda)$

select new “population” of size μ (e.g., best offspring)

update mean vector m

update global step size σ

update covariance matrix C

until stopping criterion met

Evolution Strategies

input $m \in \mathbb{R}^d$, $\sigma > 0$, $C(= I) \in \mathbb{R}^{d \times d}$

loop

sample “offspring” $x_1, \dots, x_\lambda \sim \mathcal{N}(m, \sigma^2 C)$

evaluate $f(x_1), \dots, f(x_\lambda)$

select new “population” of size μ (e.g., best offspring)

update mean vector m

update global step size σ

update covariance matrix C

until stopping criterion met

- randomized
- population-based
- rank-based (function-value free)
- step size control
- metric learning (covariance matrix adaptation, CMA)

- Why rely on a population?

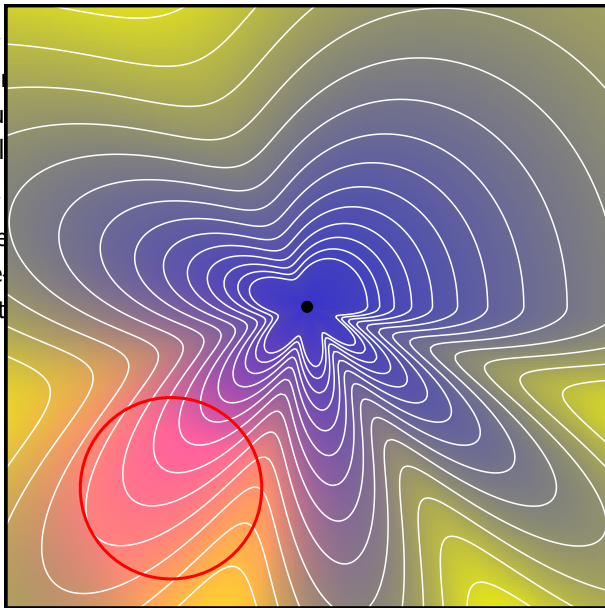
- Why rely on a population?
Not necessary, “hillclimber” $(1+1)$ -ES also works, but populations are more robust for noisy and multi-modal problems.

- Why rely on a population?
Not necessary, “hillclimber” (1+1)-ES also works, but populations are more robust for noisy and multi-modal problems.
- Why online adaptation of the step size σ ?

- Why rely on a population?
Not necessary, “hillclimber” (1+1)-ES also works, but populations are more robust for noisy and multi-modal problems.
- Why online adaptation of the step size σ ?
Scale-invariant algorithm \Rightarrow linear convergence on scale-invariant problems \Rightarrow locally linear convergence on \mathcal{C}^2 functions.

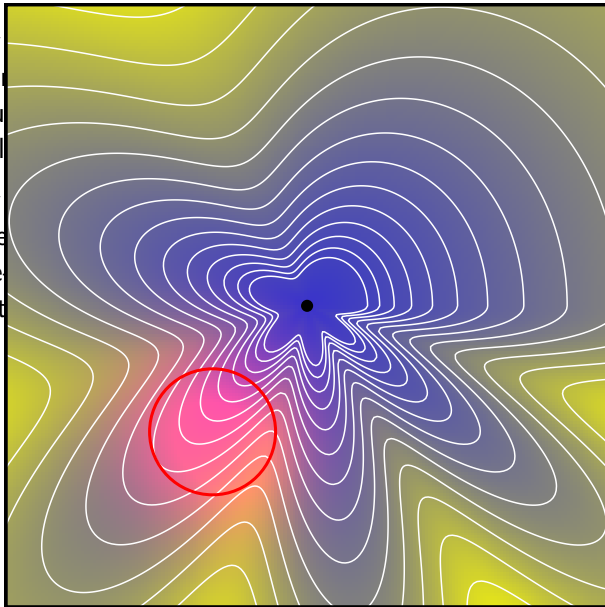
Evolution Strategies: Design Principles

- Why Not popular problem
- Why Scale scale function



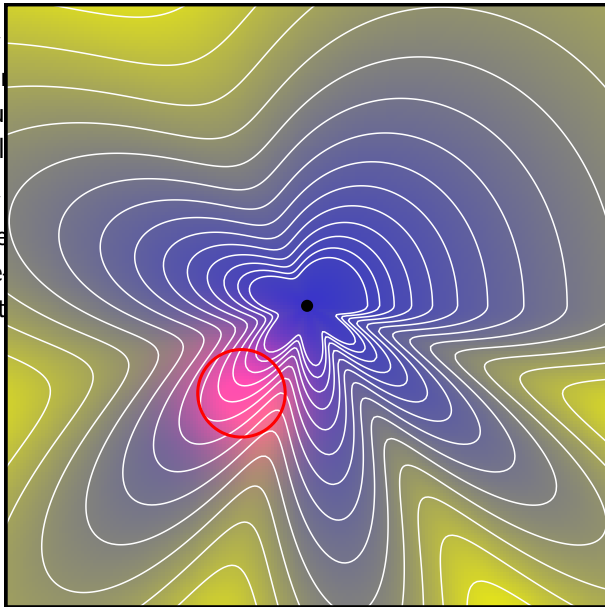
Evolution Strategies: Design Principles

- Why Not a popular problem
- Why Scale scale function



Evolution Strategies: Design Principles

- Why Not a popular problem
- Why Scale scale function

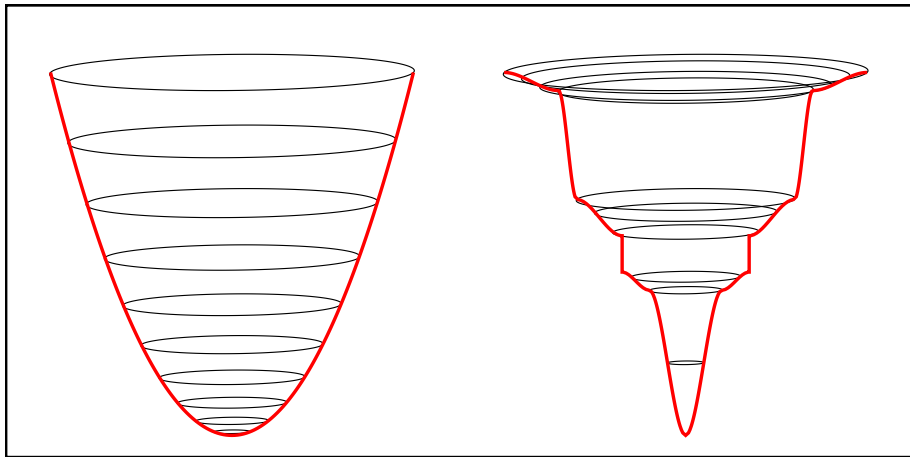


on \mathcal{C}^2

- Why rely on a population?
Not necessary, “hillclimber” (1+1)-ES also works, but populations are more robust for noisy and multi-modal problems.
- Why online adaptation of the step size σ ?
Scale-invariant algorithm \Rightarrow linear convergence on scale-invariant problems \Rightarrow locally linear convergence on \mathcal{C}^2 functions.
- Why discard function values and keep only ranks?

- Why rely on a population?
Not necessary, “hillclimber” (1+1)-ES also works, but populations are more robust for noisy and multi-modal problems.
- Why online adaptation of the step size σ ?
Scale-invariant algorithm \Rightarrow linear convergence on scale-invariant problems \Rightarrow locally linear convergence on \mathcal{C}^2 functions.
- Why discard function values and keep only ranks?
Invariance to strictly monotonic (rank-preserving) transformations of objective values.

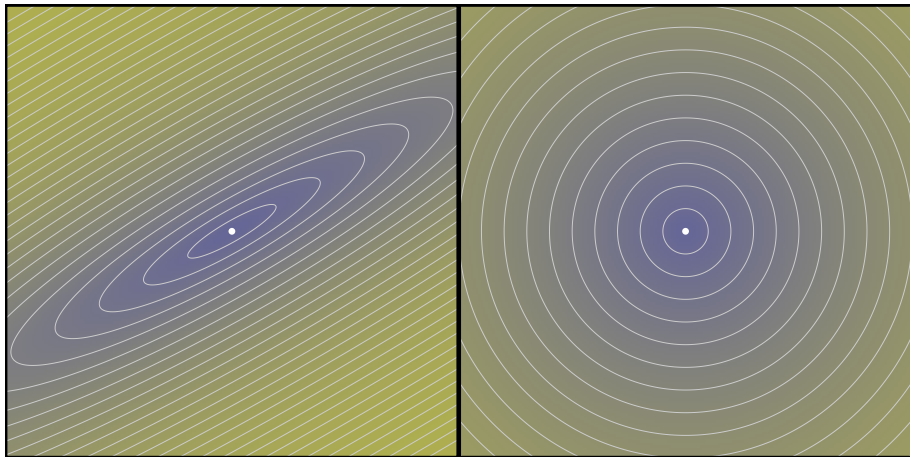
- Why rely on a population?



- Why rely on a population?
Not necessary, “hillclimber” (1+1)-ES also works, but populations are more robust for noisy and multi-modal problems.
- Why online adaptation of the step size σ ?
Scale-invariant algorithm \Rightarrow linear convergence on scale-invariant problems \Rightarrow locally linear convergence on \mathcal{C}^2 functions.
- Why discard function values and keep only ranks?
Invariance to strictly monotonic (rank-preserving) transformations of objective values.
- Why CMA?

- Why rely on a population?
Not necessary, “hillclimber” (1+1)-ES also works, but populations are more robust for noisy and multi-modal problems.
- Why online adaptation of the step size σ ?
Scale-invariant algorithm \Rightarrow linear convergence on scale-invariant problems \Rightarrow locally linear convergence on \mathcal{C}^2 functions.
- Why discard function values and keep only ranks?
Invariance to strictly monotonic (rank-preserving) transformations of objective values.
- Why CMA?
Efficient optimization of ill-conditioned problems, similar to (quasi) Newton methods.

- Why rely on a population?



Efficient optimization of ill-conditioned problems, similar to (quasi) Newton methods.

**Given samples and their ranks,
how to update the distribution parameters?**

Information Geometric Perspective

- Change of perspective: algorithm state **population** $x_1, \dots, x_\mu \rightarrow$ **distribution** $\mathcal{N}(m, C)$.

- Change of perspective: algorithm state **population** $x_1, \dots, x_\mu \rightarrow$ **distribution** $\mathcal{N}(m, C)$.
- Parameters $(m, C) = \theta \in \Theta$, distribution P_θ , pdf p_θ .

- Change of perspective: algorithm state **population** $x_1, \dots, x_\mu \rightarrow$ **distribution** $\mathcal{N}(m, C)$.
- Parameters $(m, C) = \theta \in \Theta$, distribution P_θ , pdf p_θ .
- For multi-variate Gaussians:

$$\Theta = \mathbb{R}^d \times \mathcal{P}_d$$
$$\mathcal{P}_d = \left\{ M \in \mathbb{R}^{d \times d} \mid M = M^T, M \text{ pos. def.} \right\}$$

- Change of perspective: algorithm state **population** $x_1, \dots, x_\mu \rightarrow$ **distribution** $\mathcal{N}(m, C)$.
- Parameters $(m, C) = \theta \in \Theta$, distribution P_θ , pdf p_θ .
- For multi-variate Gaussians:

$$\Theta = \mathbb{R}^d \times \mathcal{P}_d$$
$$\mathcal{P}_d = \left\{ M \in \mathbb{R}^{d \times d} \mid M = M^T, M \text{ pos. def.} \right\}$$

- **Statistical manifold** of distributions

$$\left\{ P_\theta \mid \theta \in \Theta \right\} \cong \Theta .$$

- Change of perspective: algorithm state **population** $x_1, \dots, x_\mu \rightarrow$ **distribution** $\mathcal{N}(m, C)$.
- Parameters $(m, C) = \theta \in \Theta$, distribution P_θ , pdf p_θ .
- For multi-variate Gaussians:

$$\Theta = \mathbb{R}^d \times \mathcal{P}_d$$
$$\mathcal{P}_d = \left\{ M \in \mathbb{R}^{d \times d} \mid M = M^T, M \text{ pos. def.} \right\}$$

- **Statistical manifold** of distributions

$$\left\{ P_\theta \mid \theta \in \Theta \right\} \cong \Theta .$$

- Equipped with intrinsic (Riemannian) geometry; metric tensor given by Fisher information matrix $\mathcal{I}(\theta)$.

- Goal: optimization of $\theta \in \Theta$ instead of $x \in \mathbb{R}^d$.

- Goal: optimization of $\theta \in \Theta$ instead of $x \in \mathbb{R}^d$.
- Lift objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ to objective $W_f : \Theta \rightarrow \mathbb{R}$.

- Goal: optimization of $\theta \in \Theta$ instead of $x \in \mathbb{R}^d$.
- Lift objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ to objective $W_f : \Theta \rightarrow \mathbb{R}$.
- Simplest choice:

$$W_f(\theta) = \mathbb{E}_{x \sim P_\theta} [f(x)]$$

- Goal: optimization of $\theta \in \Theta$ instead of $x \in \mathbb{R}^d$.
- Lift objective function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ to objective $W_f : \Theta \rightarrow \mathbb{R}$.
- Simplest choice:

$$W_f(\theta) = \mathbb{E}_{x \sim P_\theta} [f(x)]$$

- More flexible choice:

$$W_f(\theta) = \mathbb{E}_{x \sim P_\theta} [w(f(x))]$$

with monotonic weight function $w : \mathbb{R} \rightarrow \mathbb{R}$.

- Expectation operator “adds one degree of smoothness”.

- Expectation operator “adds one degree of smoothness”.
- Hence under weak assumptions $W_f(\theta)$ can be optimized with **gradient descent (GD)**:

$$\theta \leftarrow \theta - \eta \cdot \nabla_{\theta} W_f(\theta)$$

- Expectation operator “adds one degree of smoothness”.
- Hence under weak assumptions $W_f(\theta)$ can be optimized with **gradient descent (GD)**:

$$\theta \leftarrow \theta - \eta \cdot \nabla_{\theta} W_f(\theta)$$

- “log-likelihood trick”:

$$\begin{aligned}\nabla_{\theta} W_f(\theta) &= \nabla_{\theta} \mathbb{E}_{x \sim P_{\theta}} [w(f(x))] \\ &= \mathbb{E}_{x \sim P_{\theta}} [\nabla_{\theta} \log(p_{\theta}(x)) \cdot w(f(x))]\end{aligned}$$

- Problem: this does not work well. Why?

- Problem: this does not work well. Why?
- We have to replace the “plain” gradient on (Euclidean) parameter space Θ with the **natural gradient** respecting the intrinsic geometry of the statistical manifold:

$$\tilde{\nabla}_{\theta} W_f(\theta) = \left(\mathcal{I}(\theta) \right)^{-1} \cdot \nabla_{\theta} W_f(\theta)$$

- Problem: this does not work well. Why?
- We have to replace the “plain” gradient on (Euclidean) parameter space Θ with the **natural gradient** respecting the intrinsic geometry of the statistical manifold:

$$\tilde{\nabla}_{\theta} W_f(\theta) = \left(\mathcal{I}(\theta) \right)^{-1} \cdot \nabla_{\theta} W_f(\theta)$$

- The natural gradient is invariant under changes of the parameterization $\theta \mapsto P_{\theta}$.

- Problem: this does not work well. Why?
- We have to replace the “plain” gradient on (Euclidean) parameter space Θ with the **natural gradient** respecting the intrinsic geometry of the statistical manifold:

$$\tilde{\nabla}_{\theta} W_f(\theta) = \left(\mathcal{I}(\theta) \right)^{-1} \cdot \nabla_{\theta} W_f(\theta)$$

- The natural gradient is invariant under changes of the parameterization $\theta \mapsto P_{\theta}$.
- (Natural) gradient vector field $\Theta \rightarrow T\Theta$ defines (natural) gradient flow $\phi^t(\theta)$ tangential to $\tilde{\nabla}_{\theta} W_f(\theta)$.

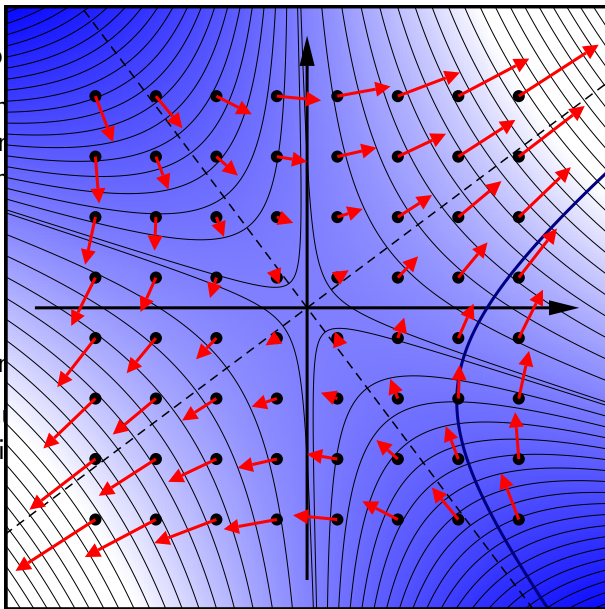
Information Geometric Perspective

- Prob

- We h
par
intrin

- The
par

- (Nat
gradi



)
ng the

al)

- Problem: this does not work well. Why?
- We have to replace the “plain” gradient on (Euclidean) parameter space Θ with the **natural gradient** respecting the intrinsic geometry of the statistical manifold:

$$\tilde{\nabla}_{\theta} W_f(\theta) = \left(\mathcal{I}(\theta) \right)^{-1} \cdot \nabla_{\theta} W_f(\theta)$$

- The natural gradient is invariant under changes of the parameterization $\theta \mapsto P_{\theta}$.
- (Natural) gradient vector field $\Theta \rightarrow T\Theta$ defines (natural) gradient flow $\phi^t(\theta)$ tangential to $\tilde{\nabla}_{\theta} W_f(\theta)$.
- Optimization: follow inverse flow curves $t \mapsto \phi^{-t}(\theta)$ from θ into (local) minimum of W_f .

- Black-box setting: expectation

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log(p_\theta(x)) \cdot w(f(x))]$$

is intractable.

- Black-box setting: expectation

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log(p_\theta(x)) \cdot w(f(x))]$$

is intractable.

- **Monte Carlo (MC)** approximation

$$\nabla_\theta W_f(\theta) \approx G(\theta) = \frac{1}{\lambda} \sum_{i=1}^{\lambda} \left[\nabla_\theta \log(p_\theta(x_i)) \cdot w(f(x_i)) \right]$$

for $x_1, \dots, x_\lambda \sim P_\theta$.

- Black-box setting: expectation

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log(p_\theta(x)) \cdot w(f(x))]$$

is intractable.

- **Monte Carlo (MC)** approximation

$$\nabla_\theta W_f(\theta) \approx G(\theta) = \frac{1}{\lambda} \sum_{i=1}^{\lambda} \left[\nabla_\theta \log(p_\theta(x_i)) \cdot w(f(x_i)) \right]$$

for $x_1, \dots, x_\lambda \sim P_\theta$.

- The estimate $\mathbb{E}[G(\theta)] = \nabla_\theta W_f(\theta)$ is consistent, hence following $G(\theta)$ amounts to **stochastic gradient descent**.

- Black-box setting: expectation

$$\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log(p_\theta(x)) \cdot w(f(x))]$$

is intractable.

- **Monte Carlo (MC)** approximation

$$\nabla_\theta W_f(\theta) \approx G(\theta) = \frac{1}{\lambda} \sum_{i=1}^{\lambda} \left[\nabla_\theta \log(p_\theta(x_i)) \cdot w(f(x_i)) \right]$$

for $x_1, \dots, x_\lambda \sim P_\theta$.

- The estimate $\mathbb{E}[G(\theta)] = \nabla_\theta W_f(\theta)$ is consistent, hence following $G(\theta)$ amounts to **stochastic gradient descent**.
- Yields natural gradient MC approximation

$$\tilde{G}(\theta) = (\mathcal{I}(\theta))^{-1} \cdot G(\theta) .$$

- Stochastic Natural Gradient Descent (SNGD) update rule:

$$\theta \leftarrow \theta - \eta \cdot \tilde{G}(\theta)$$

- Stochastic Natural Gradient Descent (SNGD) update rule:

$$\theta \leftarrow \theta - \eta \cdot \tilde{G}(\theta)$$

- SNGD is a two-fold approximation of the gradient flow:
 - discretized time: Euler steps,
 - randomized gradient: MC sampling.

- Stochastic Natural Gradient Descent (SNGD) update rule:

$$\theta \leftarrow \theta - \eta \cdot \tilde{G}(\theta)$$

- SNGD is a two-fold approximation of the gradient flow:
 - discretized time: Euler steps,
 - randomized gradient: MC sampling.
- Natural gradient flow is invariant under choice of distribution parameters $\theta \mapsto P_\theta$.

- Stochastic Natural Gradient Descent (SNGD) update rule:

$$\theta \leftarrow \theta - \eta \cdot \tilde{G}(\theta)$$

- SNGD is a two-fold approximation of the gradient flow:
 - discretized time: Euler steps,
 - randomized gradient: MC sampling.
- Natural gradient flow is invariant under choice of distribution parameters $\theta \mapsto P_\theta$.
- SNGD algorithm invariant in first order approximation due to Euler steps.

- Stochastic Natural Gradient Descent (SNGD) update rule:

$$\theta \leftarrow \theta - \eta \cdot \tilde{G}(\theta)$$

- SNGD is a two-fold approximation of the gradient flow:
 - discretized time: Euler steps,
 - randomized gradient: MC sampling.
- Natural gradient flow is invariant under choice of distribution parameters $\theta \mapsto P_\theta$.
- SNGD algorithm invariant in first order approximation due to Euler steps.

Ollivier et al. Information-Geometric Optimization Algorithms: A Unifying Picture via Invariance Principles. *arXiv:1106.3708*, 2011.

- Natural (gradient) Evolution Strategies (NES) closely follow this scheme. They were the first ES derived this way.

- Natural (gradient) Evolution Strategies (NES) closely follow this scheme. They were the first ES derived this way.
- NES approach: optimization of expected objective value with Gaussian distributions.

- Natural (gradient) Evolution Strategies (NES) closely follow this scheme. They were the first ES derived this way.
- NES approach: optimization of expected objective value with Gaussian distributions.
- Offspring x_1, \dots, x_λ act as MC sample for the estimation of $G(\theta)$.

- Natural (gradient) Evolution Strategies (NES) closely follow this scheme. They were the first ES derived this way.
- NES approach: optimization of expected objective value with Gaussian distributions.
- Offspring x_1, \dots, x_λ act as MC sample for the estimation of $G(\theta)$.
- Closed form Fisher tensor for $\mathcal{N}(m, C)$:

$$\mathcal{I}_{i,j}(\theta) = \frac{\partial m^T}{\partial \theta_i} C^{-1} \frac{\partial m}{\partial \theta_j} + \frac{1}{2} \text{tr} \left(C^{-1} \frac{\partial C}{\partial \theta_i} C^{-1} \frac{\partial C}{\partial \theta_j} \right)$$

NES algorithm

input $\theta \in \Theta$, $\lambda \in \mathbb{N}$, $\eta > 0$

loop

sample $x_1, \dots, x_\lambda \sim P_\theta$

evaluate $f(x_1), \dots, f(x_\lambda)$

$$G(\theta) \leftarrow \frac{1}{\lambda} \sum_{i=1}^{\lambda} \nabla_{\theta} \log(p_{\theta}(x_i)) \cdot f(x_i)$$

$$\tilde{G}(\theta) \leftarrow (\mathcal{I}(\theta))^{-1} \cdot G(\theta)$$

$$\theta \leftarrow \theta - \eta \cdot \tilde{G}(\theta)$$

until stopping criterion met

Wierstra et al. Natural Evolution Strategies. *CEC*, 2008 and *JMLR*, 2014.

- NES works better when replacing $f(x_i)$ in

$$G(\theta) = \frac{1}{\lambda} \sum_{i=1}^{\lambda} [\nabla_{\theta} \log(p_{\theta}(x_i)) \cdot f(x_i)]$$

with rank-based utility values w_1, \dots, w_{λ} .

- NES works better when replacing $f(x_i)$ in

$$G(\theta) = \frac{1}{\lambda} \sum_{i=1}^{\lambda} [\nabla_{\theta} \log(p_{\theta}(x_i)) \cdot f(x_i)]$$

with rank-based utility values w_1, \dots, w_{λ} .

- Sample ranks are f -quantile estimators, hence utility values can be represented as $w(f(x_i))$ with special weight function $w = w_{\theta}$ based on f -quantiles under **current** distribution P_{θ} .

- NES works better when replacing $f(x_i)$ in

$$G(\theta) = \frac{1}{\lambda} \sum_{i=1}^{\lambda} [\nabla_{\theta} \log(p_{\theta}(x_i)) \cdot f(x_i)]$$

with rank-based utility values w_1, \dots, w_{λ} .

- Sample ranks are f -quantile estimators, hence utility values can be represented as $w(f(x_i))$ with special weight function $w = w_{\theta}$ based on f -quantiles under **current** distribution P_{θ} .
- This turns NES into a function value free algorithm.

- NES works better when replacing $f(x_i)$ in

$$G(\theta) = \frac{1}{\lambda} \sum_{i=1}^{\lambda} [\nabla_{\theta} \log(p_{\theta}(x_i)) \cdot f(x_i)]$$

with rank-based utility values w_1, \dots, w_{λ} .

- Sample ranks are f -quantile estimators, hence utility values can be represented as $w(f(x_i))$ with special weight function $w = w_{\theta}$ based on f -quantiles under **current** distribution P_{θ} .
- This turns NES into a function value free algorithm.
- Benefits:
 - invariance under monotonic transformations of objective values
 - linear convergence on scale invariant problems.

xNES (exponential NES) algorithm

input $(m, A) \in \Theta$, $\lambda \in \mathbb{N}$, $\eta_m, \eta_A > 0$

loop

sample $z_1, \dots, z_\lambda \sim \mathcal{N}(0, I)$

transform $x_k, \dots, x_\lambda \leftarrow Az_k + m$

evaluate $f(x_1), \dots, f(x_\lambda)$

$\tilde{G}_m(\theta) \leftarrow \frac{1}{\lambda} \sum_{i=1}^{\lambda} w(f(x_i)) \cdot z_i$

$\tilde{G}_C(\theta) \leftarrow \frac{1}{\lambda} \sum_{i=1}^{\lambda} w(f(x_i)) \cdot \frac{1}{2}(z_i z_i^T - I)$

$m \leftarrow m - \eta_m \cdot A \cdot \tilde{G}_m(\theta)$

$A \leftarrow A \cdot \exp\left(-\eta_A \cdot \frac{1}{2} \tilde{G}_C(\theta)\right)$

until stopping criterion met

Glasmachers et al. Exponential Natural Evolution Strategies. *GECCO*, 2010.

- The resulting algorithm is indeed an ES (\mathbb{R}^d perspective), and at the same time an SNGD algorithm (Θ perspective).

- The resulting algorithm is indeed an ES (\mathbb{R}^d perspective), and at the same time an SNGD algorithm (Θ perspective).
- ES traditionally have three distinct and often very different mechanisms for
 - optimization: adaptation of m ,
 - step size control: adaptation of σ ,
 - metric learning: adaptation of C .

- The resulting algorithm is indeed an ES (\mathbb{R}^d perspective), and at the same time an SNGD algorithm (Θ perspective).
- ES traditionally have three distinct and often very different mechanisms for
 - optimization: adaptation of m ,
 - step size control: adaptation of σ ,
 - metric learning: adaptation of C .
- SNGD has only a single mechanism.

- The resulting algorithm is indeed an ES (\mathbb{R}^d perspective), and at the same time an SNGD algorithm (Θ perspective).
- ES traditionally have three distinct and often very different mechanisms for
 - optimization: adaptation of m ,
 - step size control: adaptation of σ ,
 - metric learning: adaptation of C .
- SNGD has only a single mechanism.
- Astonishing insight: all (most) parameters can be updated with a single mechanism.

- Although derived completely differently, this algorithm turns out to be closely related to CMA-ES.

Akimoto et al. Bidirectional relation between CMA evolution strategies and natural evolution strategies. PPSN 2010.

- Although derived completely differently, this algorithm turns out to be closely related to CMA-ES.

Akimoto et al. Bidirectional relation between CMA evolution strategies and natural evolution strategies. PPSN 2010.

- In particular, the update of m is identical and the rank- μ update of CMA-ES coincides with a NES update.

- Although derived completely differently, this algorithm turns out to be closely related to CMA-ES.

Akimoto et al. Bidirectional relation between CMA evolution strategies and natural evolution strategies. PPSN 2010.

- In particular, the update of m is identical and the rank- μ update of CMA-ES coincides with a NES update.
- This is astonishing, it is surely not by coincidence.

- Although derived completely differently, this algorithm turns out to be closely related to CMA-ES.

Akimoto et al. Bidirectional relation between CMA evolution strategies and natural evolution strategies. PPSN 2010.

- In particular, the update of m is identical and the rank- μ update of CMA-ES coincides with a NES update.
- This is astonishing, it is surely not by coincidence.
- This is insightful: it means that CMA-ES is (essentially) a SNGD algorithm.

Summary

- Evolution Strategies (ES) are randomized direct search methods suitable for continuous black box optimization.

- Evolution Strategies (ES) are randomized direct search methods suitable for continuous black box optimization.
- Here we have focused on a core question: how to update the distribution parameters in a principled way?

- Evolution Strategies (ES) are randomized direct search methods suitable for continuous black box optimization.
- Here we have focused on a core question: how to update the distribution parameters in a principled way?
- The parameter space is equipped with the non-Euclidean information geometry of the corresponding statistical manifold of search distributions.

- Evolution Strategies (ES) are randomized direct search methods suitable for continuous black box optimization.
- Here we have focused on a core question: how to update the distribution parameters in a principled way?
- The parameter space is equipped with the non-Euclidean information geometry of the corresponding statistical manifold of search distributions.
- SNGD on a stochastically relaxed problem results in a direct search algorithm: the Natural-gradient Evolution Strategy (NES) algorithm.

- Evolution Strategies (ES) are randomized direct search methods suitable for continuous black box optimization.
- Here we have focused on a core question: how to update the distribution parameters in a principled way?
- The parameter space is equipped with the non-Euclidean information geometry of the corresponding statistical manifold of search distributions.
- SNGD on a stochastically relaxed problem results in a direct search algorithm: the Natural-gradient Evolution Strategy (NES) algorithm.
- The SNGD parameter update is by no means restricted to Gaussian distributions. It is a general construction template for update equations of continuous distribution parameters.

Thank you!
Questions?