

University of Paris Sud 11, France
Department of computer science
Machine learning, information and contents (AIC)
Stochastic optimization class

Theoretical foundation for covariance matrix adaptation evolutionary strategy (CMA-ES) from information geometry perspective

Presented by : Ahmed MAZARI

20 February 2017

Context and motivation

Context :

Continuous domain optimization

Target :

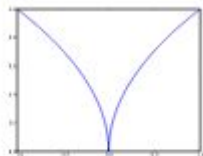
Minimize an objective function in a continuous domain

$$f: X \subseteq \mathbb{R}^n \longrightarrow \mathbb{R}, \quad x \longmapsto f(x)$$

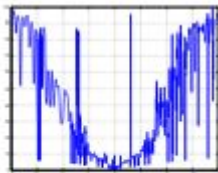
Context and motivation

Black box scenario :

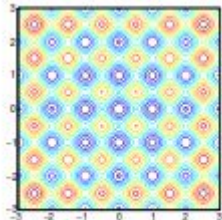
How to minimize f in a black-box scenario ?
Why it's difficult to solve ?



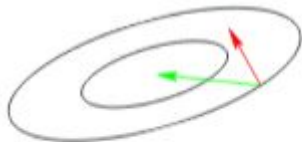
non-convex



ruggedness



non-separability



gradient direction Newton direction

ill-conditioning



- Gradient not available
- Non-convex, noisy, rugged, high dimensional
- Analytic form is not known
- Time evaluation ...

Context and motivation

Goal :

Convergence to the global optimum.

As fast as possible

Problem :

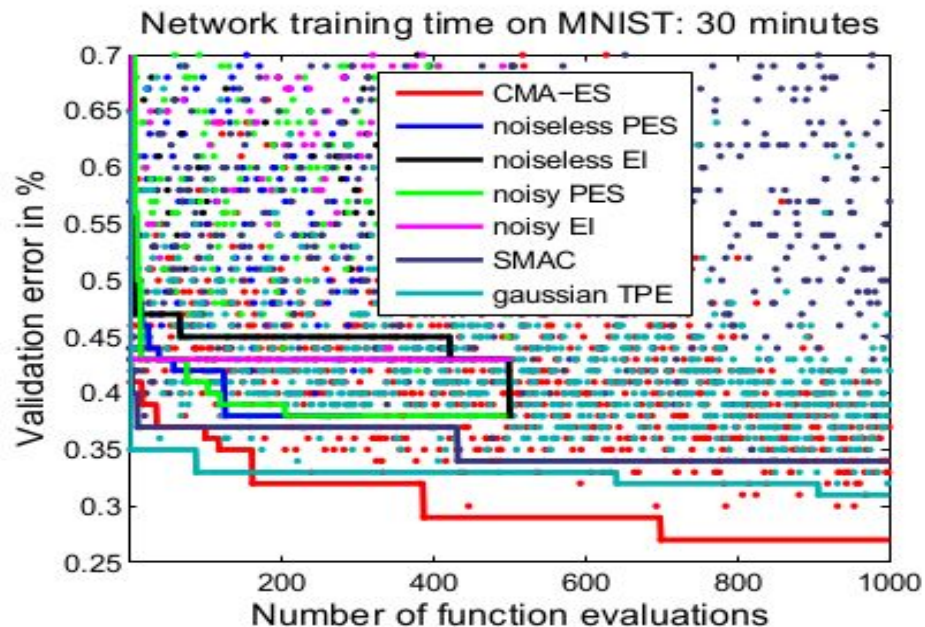
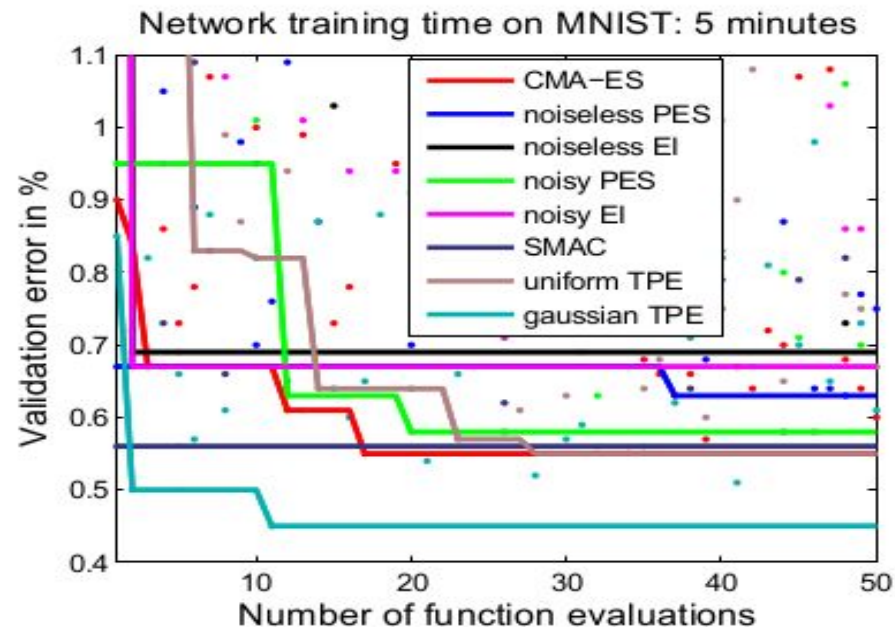
- Exhaustive search is infeasible
- Naive random and deterministic search take too long time

CMA-ES for deep learning hyperparameter optimization

name	description	transformation	range
x_1	selection pressure at e_0	$10^{-2+10^2x_1}$	$[10^{-2}, 10^{98}]$
x_2	selection pressure at e_{end}	$10^{-2+10^2x_2}$	$[10^{-2}, 10^{98}]$
x_3	batch size at e_0	2^{4+4x_3}	$[2^4, 2^8]$
x_4	batch size at e_{end}	2^{4+4x_4}	$[2^4, 2^8]$
x_5	frequency of loss recomputation r_{freq}	$2x_5$	$[0, 2]$
x_6	alpha for batch normalization	$0.01 + 0.2x_6$	$[0.01, 0.21]$
x_7	epsilon for batch normalization	10^{-8+5x_7}	$[10^{-8}, 10^{-3}]$
x_8	dropout rate after the first Max-Pooling layer	$0.8x_8$	$[0, 0.8]$
x_9	dropout rate after the second Max-Pooling layer	$0.8x_9$	$[0, 0.8]$
x_{10}	dropout rate before the output layer	$0.8x_{10}$	$[0, 0.8]$
x_{11}	number of filters in the first convolution layer	$2^{3+5x_{11}}$	$[2^3, 2^8]$
x_{12}	number of filters in the second convolution layer	$2^{3+5x_{12}}$	$[2^3, 2^8]$
x_{13}	number of units in the fully-connected layer	$2^{4+5x_{13}}$	$[2^4, 2^9]$
x_{14}	Adadelata: learning rate at e_0	$10^{0.5-2x_{14}}$	$[10^{-1.5}, 10^{0.5}]$
x_{15}	Adadelata: learning rate at e_{end}	$10^{0.5-2x_{15}}$	$[10^{-1.5}, 10^{0.5}]$
x_{16}	Adadelata: ρ	$0.8 + 0.199x_{16}$	$[0.8, 0.999]$
x_{17}	Adadelata: ϵ	$10^{-3-6x_{17}}$	$[10^{-9}, 10^{-3}]$
x_{14}	Adam: learning rate at e_0	$10^{-1-3x_{14}}$	$[10^{-4}, 10^{-1}]$
x_{15}	Adam: learning rate at e_{end}	$10^{-3-3x_{15}}$	$[10^{-6}, 10^{-3}]$
x_{16}	Adam: β_1	$0.8 + 0.199x_{16}$	$[0.8, 0.999]$
x_{17}	Adam: ϵ	$10^{-3-6x_{17}}$	$[10^{-9}, 10^{-3}]$
x_{18}	Adam: β_2	$1 - 10^{-2-2x_{18}}$	$[0.99, 0.9999]$
x_{19}	adaptation end epoch index e_{end}	$20 + 200x_{19}$	$[20, 220]$

Hyperparameters to optimize

CMA-ES for deep learning hyperparameter optimization



Comparison of optimizers for Adam with batch selection when solutions are evaluated sequentially for 5 minutes each (left), and in parallel for 30 minutes each (right).

Covariance matrix adaptation evolutionary strategy (CMA-ES) algorithm

The CMA-ES

Input: $\mathbf{m} \in \mathbb{R}^n$, $\sigma \in \mathbb{R}_+$, λ

Initialize: $\mathbf{C} = \mathbf{I}$, and $\mathbf{p}_c = \mathbf{0}$, $\mathbf{p}_\sigma = \mathbf{0}$,

Set: $c_c \approx 4/n$, $c_\sigma \approx 4/n$, $c_1 \approx 2/n^2$, $c_\mu \approx \mu_w/n^2$, $c_1 + c_\mu \leq 1$, $d_\sigma \approx 1 + \sqrt{\frac{\mu_w}{n}}$,
and $w_{i=1 \dots \lambda}$ such that $\mu_w = \frac{1}{\sum_{i=1}^{\lambda} w_i^2} \approx 0.3 \lambda$

While not terminate

$\mathbf{x}_i = \mathbf{m} + \sigma \mathbf{y}_i$, $\mathbf{y}_i \sim \mathcal{N}_i(\mathbf{0}, \mathbf{C})$, for $i = 1, \dots, \lambda$ sampling

$\mathbf{m} \leftarrow \sum_{i=1}^{\mu} w_i \mathbf{x}_{i:\lambda} = \mathbf{m} + \sigma \mathbf{y}_w$ where $\mathbf{y}_w = \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda}$ update mean

$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \mathbb{1}_{\{\|\mathbf{p}_\sigma\| < 1.5\sqrt{n}\}} \sqrt{1 - (1 - c_c)^2} \sqrt{\mu_w} \mathbf{y}_w$ cumulation for \mathbf{C}

$\mathbf{p}_\sigma \leftarrow (1 - c_\sigma) \mathbf{p}_\sigma + \sqrt{1 - (1 - c_\sigma)^2} \sqrt{\mu_w} \mathbf{C}^{-\frac{1}{2}} \mathbf{y}_w$ cumulation for σ

$\mathbf{C} \leftarrow (1 - c_1 - c_\mu) \mathbf{C} + c_1 \mathbf{p}_c \mathbf{p}_c^T + c_\mu \sum_{i=1}^{\mu} w_i \mathbf{y}_{i:\lambda} \mathbf{y}_{i:\lambda}^T$ update \mathbf{C}

$\sigma \leftarrow \sigma \times \exp\left(\frac{c_\sigma}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1\right)\right)$ update of σ

Not covered on this slide: termination, restarts, useful output, boundaries and encoding

Rank-u update in CMA-ES

- Increases the possible learning rate in large populations
- Uses the evolution path and reduces the number of necessary function evaluations

Properties of CMA-ES

- Learns the dependencies between variables
- Invariant for any increasing function $g : g(f(x))$
- increases the likelihood of previously successful steps
- Approximates the inverse of hessian on quadratic functions

Problems of CMA-ES

- $O(n^2)$ time and space complexities :
 - To store and update C in $R^{(n \text{ by } n)}$
 - To compute the eigendecomposition of C
- Internal CPU-time $10^{(-8)} n^2$ secondes per function evaluation on a 2GHZ pc, tweaks are available
 - 1 000 000 f -evaluations in 100-D take 100 seconds internal CPU-time

Youhei Akimoto, Anne Auger and Nikolaus Hansen . In CMA-ES and Advanced Adaptation Mechanisms talk

Covariance matrix adaptation evolutionary strategy (CMA-ES)

Initialize distribution parameters $\{m, C\}$, set population size $\lambda \in \mathbb{N}$

While not terminate

1- **Sample** λ independent points $x_1 \dots x_\lambda$ from $P(x|\theta)$ such that $P(x|\theta)$ multivariate gaussian distribution

$\theta = \{m, \sigma^2, C\}$

m : mean vector, $m \in \mathbb{R}^n$

σ : global step size controls step length, $\sigma \in \mathbb{R}_+$

C : Covariance matrix (symmetric, positive definite) determines the shape of the distribution ellipsoid, $C \in \mathbb{R}^{n \times n}$

2- **Evaluate** the fitness values $f(x_1), \dots, f(x_\lambda)$

3- **Update** the parameters of θ

$$m^{t+1} = m^t + \eta_m \sum_{i=1}^{\lambda} W_{R_i} (x_i - m^t) \quad (1)$$

$$C^{t+1} = C^t + \eta_C \sum_{i=1}^{\lambda} W_{R_i} ((x_i - m^t)(x_i - m^t)^T - C^t) \quad (2)$$

where η_m and η_C are learning rate parameters
and W_{R_i} : the weight for the R_i^{th} highest point

Rank-u update CMA-ES and Natural Gradient Ascent

1- Expected Fitness

$$J(\theta) = \mathbb{E}[f(x); \theta] = \int f(x)P(x; \theta)dx \quad (3)$$

$J(\theta)$ a function on a Riemannian manifold.

2-Natural Gradient

$$\tilde{\nabla} J(\theta) = F_{\theta}^{-1} \nabla J(\theta) \quad (4)$$

2-A F_{θ} : Fisher metric for θ such that

$$\begin{aligned} F_{\theta} &= \int \frac{\partial \ln P(x; \theta)}{\partial \theta} \left(\frac{\partial \ln P(x; \theta)}{\partial \theta} \right)^T P(x; \theta) dx \\ &= \mathbb{E} \left[\frac{\partial \ln P(x; \theta)}{\partial \theta} \left(\frac{\partial \ln P(x; \theta)}{\partial \theta} \right)^T \right] \end{aligned} \quad (5)$$

2-B $\nabla J(\theta)$ gradient of J

$$\begin{aligned} \nabla J(\theta) &= \nabla \int f(x)P(x, \theta) \nabla \ln P(x, \theta) dx \\ &= \mathbb{E}[f(x) \nabla \ln P(x, \theta)] \end{aligned} \quad (6)$$

Under some regularity conditions which are derived from Lebsgue's dominated convergence theorem 16.3

Billingsley, Probability and Measure. 1995

Rank-u update CMA-ES and Natural Gradient Ascent

3- Monte carlo approximation of the natural gradient : fitness function is unknown

$$\delta(\theta|\{x_i\}) = \sum_{i=1}^{\lambda} \frac{f(x_i)}{\lambda} \mathbf{F}^{-1}(\theta) \nabla \ln P(x_i; \theta) \quad (7)$$

Circumvent the computation of the inverse of the information matrix by theorem 4.1

$$\tilde{\delta}(\theta|\{x_i\}) = \sum_{i=1}^{\lambda} \frac{f(x_i)}{\lambda} \left[\text{vech}((x_i - m(\theta^t))(x_i - m(\theta^t))^T - C(\theta^t)) \right] \quad (8)$$

4- Update rule for natural gradient learning :

update θ such that : $\theta = \theta + \eta \tilde{\delta}(\theta)$

$$m^{t+1} = m^t + \eta \sum_{i=1}^{\lambda} \frac{f(x_i)}{\lambda} (x_i - m^t) \quad (9)$$

$$C^{t+1} = C^t + \eta \sum_{i=1}^{\lambda} \frac{f(x_i)}{\lambda} ((x_i - m^t)(x_i - m^t)^T - C^t) \quad (10)$$

$$\eta_C = \eta_m = \eta$$

Important property of Natural Gradient

The natural gradient equipped with the Fisher metric on the density p is **invariant** under re-parameterization of the distribution

$$\begin{aligned} F_{\theta} &= \int \frac{\partial \ln P(x; \theta)}{\partial \theta} \left(\frac{\partial \ln P(x; \theta)}{\partial \theta} \right)^T P(x; \theta) dx \\ &= \mathbb{E} \left[\frac{\partial \ln P(x; \theta)}{\partial \theta} \left(\frac{\partial \ln P(x; \theta)}{\partial \theta} \right)^T \right] \end{aligned}$$

Importance of invariance

Why it is important to build invariants ?

Natural gradient ascent along with rank-u CMA-ES update

Natural gradient ascent with monte-carlo approximation \Leftrightarrow rank-u CMA-ES

From equation (1), (2) of CMA-ES

$$m^{t+1} = m^t + \eta_m \sum_{i=1}^{\lambda} W_{R_i} (x_i - m^t) \quad (1)$$

$$C^{t+1} = C^t + \eta_C \sum_{i=1}^{\lambda} W_{R_i} ((x_i - m^t)(x_i - m^t)^T - C^t) \quad (2)$$

And from (9), (10) of natural gradient

$$W_{R_i} = \sum_{i=1}^{\lambda} \frac{f(x_i)}{\lambda}$$

$$m^{t+1} = m^t + \eta_f \sum_{i=1}^{\lambda} \frac{f(x_i)}{\lambda} (x_i - m^t) \quad (9)$$

$$C^{t+1} = C^t + \eta_f \sum_{i=1}^{\lambda} \frac{f(x_i)}{\lambda} ((x_i - m^t)(x_i - m^t)^T - C^t) \quad (10)$$

Building invariants

How to build invariants ?

Building invariants

How to build invariants ?

- Riemannian geometry
- Information theory
- Group theory
- Approximation theory
- Theory of complexity

Perspectives and open problems

- Theoretical foundation for the evolution path and cumulation
- Evaluation of the estimated natural gradient (Coefficients)
- Stability of CMA-ES
- Fitness shaping exploration
- Build a group of invariants

Limit of CMA-ES



What if x is on an arbitrary space ?

Extended results from Information-Geometric Optimization and group theory

Perspective from Information geometric optimization :

- Optimization on arbitrary space
- Invariance principles for generalization



f-invariance , theta-invariance, X-invariance

- Quantile rewriting of the objective function

Yann Ollivier, Ludovic Arnold, Anne Auger and Nikolaus Hansen. Information-Geometric Optimization Algorithms: A Unifying Picture via Invariance Principles. 2011

Extended results from Information-Geometric Optimization and group theory

Perspective from group theory :

- Learning group of invariants (translation, rotation, elastic deformation)
- Linearization of symmetries in high dimension without a loss of information
- Retrieve more invariants from the inner structure of information

My intuition

- Building invariants in CMA-ES can also lead to theoretical grounding of deep neural networks
- Both CMA-ES and neural networks are black-box architectures
- Rely on neural networks to learn invariants for CMA.
- Apply CMA-ES to optimize the hyperparameters of deep architecture to learn the complex inner structure
- Dialectic relationship between CMA-ES and deep neural networks

References

- Youhei Akimoto, Yuichi Nagata, Isao Ono, Shigenobu Kobayashi. [Theoretical foundation for CMA-ES from information geometric perspective](#). 2012
- Yann Ollivier, Ludovic Arnold, Anne Auger and Nikolaus Hansen. [Information-Geometric Optimization Algorithms: A Unifying Picture via Invariance Principles](#). 2011
- Stephane Mallat. [Understanding deep convolutional neural networks](#). 2016
- Tom Schaul . [Studies in Continuous Black-box Optimization, PHD thesis](#). Technical university of Munich. 2011
- Ilya Loshchilov & Frank Hutter. [CMA-ES for hyperparameter optimization of deep neural networks](#). ICLR 2016
- Yann Ollivier. [Optimization and recurrent neural networks training talk](#). College de France, 2016
- Stephane Mallat. [Mathematical mysteries of deep neural networks](#). College de France, 2016
- Anne Auger. [How information theory sheds new light on black-box optimization talk](#). Institut Henré Poincaré, France, 2016
- Anne Auger, Youhei Akimoto and Nikolaus Hansen. Miscellaneous talks.
- Anne Auger and Dimo Brockhoff. Introduction and advanced optimization class (M2 AIC). 2016-2017