

Covariance Matrix Adaptation in Python

ME498: Comp. modeling & optimization

Tejaswin Parthasarathy, Mattia Gazzola

February 21, 2019

Matplotlib

Additional packages

- Seaborn DEMO
- bokeh
- plotly
- See Python wiki for more plotting tools

GA demo : Results of our GA implementation

How to improve convergence?

- Parameters?
- Strategies? (More mutation, less recombination say...?)
- Tuning is painful

⇒ CMAs (and other algorithms)

Implementation of CMAEs

The algorithm

Recap::

$$\mathbf{x}_i = m + \sigma \mathbf{z}_i$$

$$\mathbf{m} \leftarrow \mathbf{m} + \sigma \langle \mathbf{z} \rangle_w$$

$$\mathbf{p}_c \leftarrow (1 - c_c) \mathbf{p}_c + \sqrt{1 - (1 - c_c)^2} \sqrt{\frac{1}{\sum_{i=1}^{\mu} w_i^2}} \langle \mathbf{z} \rangle_w$$

$$\mathbf{C} \leftarrow (1 - c_{cov}) \mathbf{C} + c_{cov} \frac{1}{\mu_{cov}} \mathbf{p}_c \mathbf{p}_c^T + c_{cov} \left(1 - \frac{1}{\mu_{cov}} \right) \mathbf{Z}$$

$$\mathbf{p}_\sigma \leftarrow (1 - c_\sigma) \mathbf{p}_\sigma + \sqrt{1 - (1 - c_\sigma)^2} \sqrt{\frac{1}{\sum_{i=1}^{\mu} w_i^2}} \mathbf{C}^{-\frac{1}{2}} \langle \mathbf{z} \rangle_w$$

$$\sigma \leftarrow \sigma e^{\left(\frac{1}{d_\sigma} \left(\frac{\|\mathbf{p}_\sigma\|}{E \|\mathcal{N}(\mathbf{0}, \mathbf{I})\|} - 1 \right) \right)}$$

$$\mathbf{z}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$$

$$\langle \mathbf{z} \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}$$

$$\mathbf{Z} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda} \mathbf{z}_{i:\lambda}^T$$

$$\mu_{cov} = \sqrt{\frac{1}{\sum_{i=1}^{\mu} w_i^2}}$$

Figure 1: CMAEs

Problem independent

- Set evolution paths $\mathbf{p}_\sigma = \mathbf{0}, \mathbf{p}_c = \mathbf{0}$
- Set number of generations $g = 0$
- Set covariance matrix $\mathbf{C} = \mathbf{I}$ (Why?)

Problem dependent

- Distribution mean $\mathbf{m} \in \mathbb{R}^n$
- Step size $\sigma \in \mathbb{R}_{>0}$ (Important to set > 0)

Starting CMAEs: more on problem dependent parameters

- Optimum presumably be within the initial cube $\mathbf{m} \pm 3\sigma (1, 1, \dots, 1)^T$
- \therefore if optima $\in [a, b]^n$ choose $\mathbf{m} \in [a, b]^n$ (as a uniformly random vector) and $\sigma = 0.3 * (b - a)$
- Different search intervals Δs_i for different variables can be done using \mathbf{C} as shown below (deferred discussion):

$$\begin{bmatrix} \Delta s_1^2 & 0 & \cdots & 0 \\ 0 & \Delta s_2^2 & \cdots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & \Delta s_n^2 \end{bmatrix} \quad (1)$$

- Δs_i all must be of similar magnitude (for conditioning). Else, rescale your problem.

First step : Sampling

New population of points, for $k = 1 \dots \lambda$

- $\mathbf{y}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- $\mathbf{z}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{C}) = \mathbf{B}\mathbf{D}\mathbf{y}_k$
 - Given $\mathbf{C} = \mathbf{B}\mathbf{D}^2\mathbf{B}^T$
 - Consult ¹ for a proof of why $\mathbf{A}\mathcal{N}(\mathbf{0}, \mathbf{I}) = \mathcal{N}(\mathbf{0}, \mathbf{A}\mathbf{A}^T)$
- $\mathbf{x}_k = \mathbf{m} + \sigma\mathbf{z}_k \sim \mathcal{N}(\mathbf{m}, \sigma^2\mathbf{C})$

Computing?

- Steps 1 and 3 above using `*/np.multiply` for elementwise multiplication and `+` for elementwise addition

We need a way to sample correlated (across dimensions) populations from `numpy`_____

¹<https://math.stackexchange.com/q/2115701>

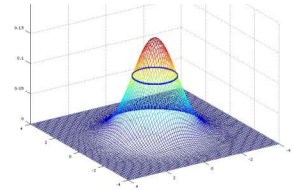
How to sample a multivariate normal distribution?

- `np.random.multivariate_normal` DEMO

Caveats?

- What is **C** / **cov** (in a 2D case) and its meaning?
- **cov** needs to be SP(semi)D. Is it? What about the update step?
- What happens in **numpy** if it is not?

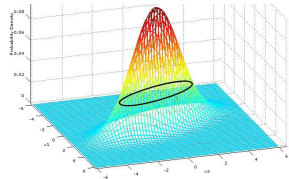
Covariance-Matrix Adaptation (CMA)



$N(0, \mathbf{I})$

Identity Matrix

How??



$N(0, \mathbf{C})$

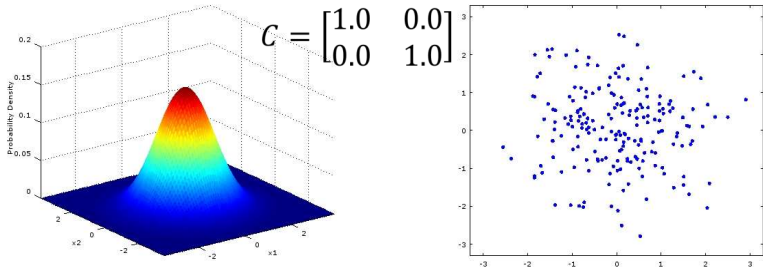
Adapt the Covariance Matrix

6/3/2016

CAIRO UNIVERSITY - COMPUTER ENGINEERING - 2015

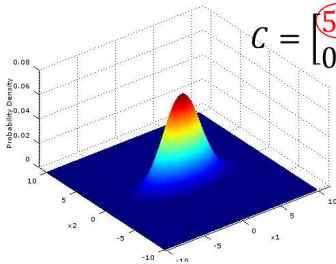
68

Covariance-Matrix

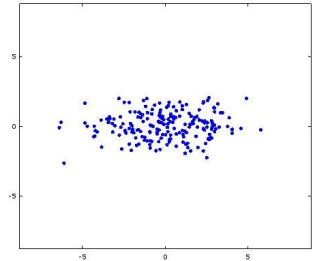


Covariance Matrix : Example 2

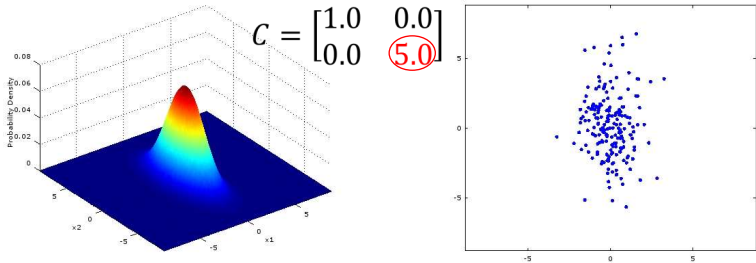
Covariance-Matrix



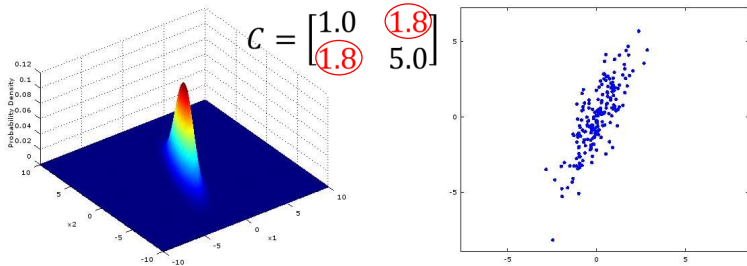
$$C = \begin{bmatrix} 5.0 & 0.0 \\ 0.0 & 1.0 \end{bmatrix}$$



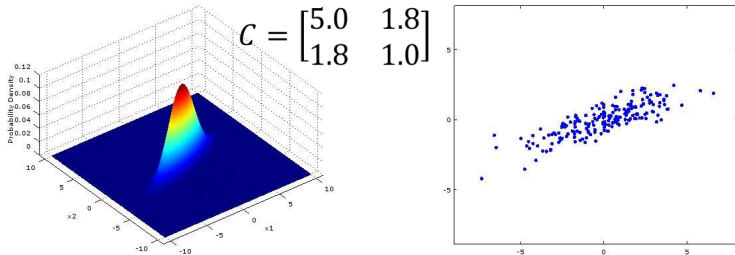
Covariance-Matrix



Covariance-Matrix

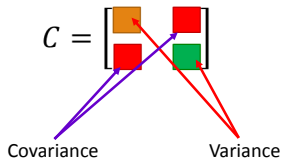


Covariance-Matrix



Covariance-Matrix Adaptation (CMA)

- To which direction should the population be directed?



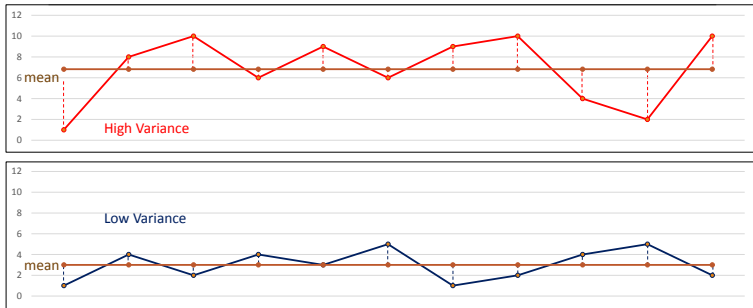
Variance

- Variance is a measure of how far a variable changes away from its mean

$$\text{var}(X) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \quad , \bar{X} \text{ is the mean of the samples of } X$$

Covariance Matrix : Example

Variance



6/3/2016

CAIRO UNIVERSITY - COMPUTER ENGINEERING - 2015

76

Covariance

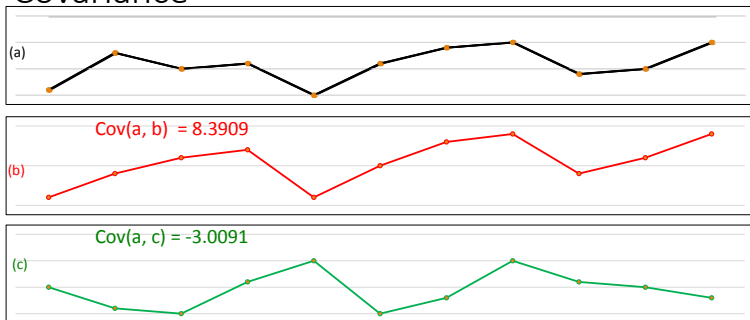
- Covariance is a measure of how two variables change together

$$\text{covar}(X, Y) = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

$$\text{covar}(X, X) = \text{var}(X)$$

Covariance Matrix : Example

Covariance



6/3/2016

CAIRO UNIVERSITY - COMPUTER ENGINEERING - 2015

78

- What is **C** / **cov** (in a 2D case) and its meaning?
 1. Covariance, how a gene varies with another (across dimensions)
 2. $\mu_x = \frac{1}{N} \sum_{i=1}^N x_i$
 3. $\sigma_x^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_x)^2$ and $\sigma_{xy} = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)$
- **cov** needs to be SPD. Is it? What about the update step?
 1. Symmetric by definition
 2. Symmetric after update too
- What happens in **numpy** if it is not?
 1. **numpy** checks for PD, else throws an exception

What is CMA-ES doing?

- How does CMA estimate \mathbf{C} ?
- What about the choice of weights?
- What is CMA doing by adapting \mathbf{C} ?

Sampling : Idea of \mathbf{C} \rightarrow math-Answers

- How does CMA estimate \mathbf{C} ?
 1. You can use the new population to get \mathbf{C} too, but information is lost (no information on how the population “evolved”, see EMNA from previous slides)
 2. **Idea** : Use $\mu_x^{(g+1)} = \frac{1}{N_{best}} \sum_{i=1}^{N_{best}} x_i$ rather than $\sigma_x^{2,(g+1)} = \frac{1}{N_{best}} \sum_{i=1}^{N_{best}} (x_i - \mu_x^{(g)})^2$, across N_{best} individuals to estimate covariance between genes (rank μ update)
 3. Exponential weighting, discussed later on
- What about the choice of weights?
 1. Reflect normalization (relates back to the ability of CMA to maintain invariance)
- What is CMA doing by adapting \mathbf{C} ?
 1. Conducts PCA (eigenvectors), rotated representation $\mathbf{C} = \mathbf{B}\mathbf{D}^2\mathbf{B}^T$, inverse Hessian (second order)

CMAEs performs PCA on the optimization data

PCA?

1. Principal Component Analysis
2. Find directions with
 - High Variance
 - Low Covariance with other components
3. Find dimensions that are “independent” from one another
4. Gives a useful basis (in this case for \mathbf{C})

Choice of λ ?

- Look at the CMA tutorial : The CMA tutorial/CMA tutorial on Arxiv
- Usually $\lambda = \lfloor 4 + 3 \ln n \rfloor$
- And $\mu = \lfloor \lambda/2 \rfloor$

Second step : Selection

How to select μ best individuals

- $\langle \mathbf{z}_k \rangle_w = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda}$
- Constraint on weights: $\sum_{i=1}^{\mu} w_i = 1, w_i > 0 \forall i = 1 \dots \mu$ (at least in our version of CMA)

Computing?

- Fitness function evaluation left upto user (including constraints etc.). This determines the μ best individuals.
- The weighted sum can be evaluated using `np.inner()/broadcasting with */np.sum()` after `*` ...

Choice of w_i ?

- Look at the CMA tutorial : The CMA tutorial/CMA tutorial on Arxiv

Third step : Recombination

Recombination to get new m

- $m \leftarrow m + \sigma \langle z \rangle_w$
- No parameters in this step!

Notice!

- σ is the “overall” step size and is a scalar.
- It could also be a matrix. Is this a good idea?
 - What about a diagonal matrix?

Computing?

- Use elementwise addition using $+$ operator

Third step : Recombination–Answers

- σ could also be a matrix. Is this a good idea?
 - NO!
 - One dimension depends on another, but not during sampling. This degrades the convergence of the algorithm
- What about a diagonal matrix?
 - NO!
 - $\mathbf{C} = \mathbf{B}\mathbf{D}^2\mathbf{B}^T$ does the job of maintaining scaling, orientation etc. of the elements.

CONCLUSION—Scalar σ is apt.

Fourth step : Step size control

Control for σ and cumulation \mathbf{p}_σ

- $\mathbf{p}_\sigma \leftarrow (1 - c_\sigma)\mathbf{p}_\sigma + \sqrt{c_\sigma(2 - c_\sigma)}\mu_{\text{cov}}\mathbf{C}^{-\frac{1}{2}}\langle \mathbf{z} \rangle_w$
- $\sigma \leftarrow \sigma \exp\left(\frac{c_\sigma}{d_\sigma}\left[\frac{\|\mathbf{p}_\sigma\|}{\mathbb{E}\|\mathcal{N}(0, \mathbf{I})\|} - 1\right]\right)$

Computing/Python?

- Notice you need to invert the covariance matrix! How will you do it?
 - **Hint:** Exploit properties of \mathbf{C} !
 - This means you just need `np.linalg.eigh()` for now (there are many other powerful methods for general symmetric matrix inverse)
 - Can reduce $\mathcal{O}(n^3)$ to $\mathcal{O}(n^2)$ in practice? (See B2. Strategy internal numerical effort in CMA tutorial)

Computing continued

- Extensive use of matvecs (\mathbf{a})
- What about the norm in the σ update?
 - What is a norm?
 - So what norm should we use?
 - The two-norm is widely used (Euclidean distance)
- What's E?
 - What's $E \|\mathcal{N}(0, \mathbf{I})\|$?
 - $\approx \sqrt{n} \left(1 - \frac{1}{4n} + \frac{1}{21n^2}\right)$

What is path update doing?

- Increase probability of reproducing successful solution paths...
- Weighting with exponential decay...
- What about the choice of weights?
 - Makes the expected length independent of the direction
 - “Follows” the random choice of $p_{\sigma}^{(0)}$

What is σ update doing?

- Decrease/Increase size until path steps are uncorrelated...
- How does the two norm of the path reflect this “un”correlation?
- What about the choice of weights?

Choice of c_σ, d_σ ?

- Look at the CMA tutorial : The CMA tutorial/CMA tutorial on Arxiv
- c_σ is learning rate for cumulation usually set to $\approx \frac{4}{n}$
- d_σ is the damping parameter for step size update
 $\approx 1 + \frac{\mu_{\text{cov}}}{\sqrt{n}}$

Fifth step : Covariance matrix adaptation

Control for \mathbf{C} and cumulation p_c

- $p_c \leftarrow (1 - c_c)p_c + \sqrt{c_c(2 - c_c)}\mu_{\text{cov}}\langle \mathbf{z} \rangle_w$
- $\mathbf{C} \leftarrow (1 - c_{\text{cov}})\mathbf{C} + \frac{c_{\text{cov}}}{\mu_{\text{cov}}}\mathbf{p}_c\mathbf{p}_c^T + c_{\text{cov}}\left(1 - \frac{1}{\mu_{\text{cov}}}\right)\mathbf{Z}$ where
 $\mathbf{Z} = \sum_{i=1}^{\mu} w_i \mathbf{z}_{i:\lambda} \mathbf{z}_{i:\lambda}^T$

Computing/Python?

- Usual operations ($*$, $+$)
- For calculating outer products, use `np.outer()`

What is cumulation for p_c doing?

- Weighting with exponential decay for prior values
- New information from PCA of steps updated into \mathbf{C} path
- What about the choice of weights?

What is \mathbf{C} update doing?

- Weighting with exponential decay for prior values
- Rank one update using p_c (What's rank?)
 - Why is the update rank one? (One-dimensional information)
 - Why use p_c rather than $\langle z \rangle$?
- Rank μ update
 - As seen earlier, CMA cleverly estimates \mathbf{C} using old step information

Choice of c_C, c_{cov} ?

- Look at the CMA tutorial : The CMA tutorial/CMA tutorial on Arxiv
- c_C is learning rate for path cumulation set to $\approx \frac{4}{n}$
- $c_{cov} \approx \frac{2+\mu_{cov}^2}{n^2}$

Terminating CMA

Algorithm should be stopped when CPU-time is wasted. Then we can:

1. restart (eventually with increased population size)
2. reconsider encoding and objective function formulation

Problem independent

- **NoEffectAxis** : Stop if adding 0.1 std.dev. vector to any direction of basis **B** does not change **m**
- **NoEffectCoord** : Stop if adding 0.2 std.dev. to any coordinate does not change **m**
- **ConditionCov**: stop if condition number of covariance matrix exceeds 10^{14}
 - Whats condition number of a matrix?
 - `np.linalg.cond()`, although you can directly check **D**

Problem independent

- **EqualFunValues**: stop if the range of the best $f(\mathbf{x})$ of the last $10 + \lceil 30n/\lambda \rceil$ generations is zero.
- **Stagnation**: Track history of the best and the median fitness in each iteration over the last 20% but at least $120 + 30n/\lambda$ and no more than 20000 iterations. Stop, if in both histories the median of the last (most recent) 30% values is not better than the median of the first 30%.
- **TolXUp**: stop if $\sigma \times \max(\text{diag}(\mathbf{D}))$ increased by more than 10^4 . This indicates a far too small initial σ , or divergence.

We note that there are problem dependent diagnostics too!

Boundaries/Constraints in CMA : Best solution strictly inside

- Set fitness (for minimization problem) as

$$f_{\text{fitness}}(\mathbf{x}) = f_{\text{max}} + \|\mathbf{x} - \mathbf{x}_{\text{feasible}}\|$$

1. Notation

1.1 f_{max} is larger than worst feasible fitness

1.2 $\mathbf{x}_{\text{feasible}}$ is constant, in the middle of feasible region

2. Caveat : Optimal solution not too close to the infeasible region

- Alternatively, resample any infeasible point until it becomes feasible

Boundaries/Constraints in CMA : Repair

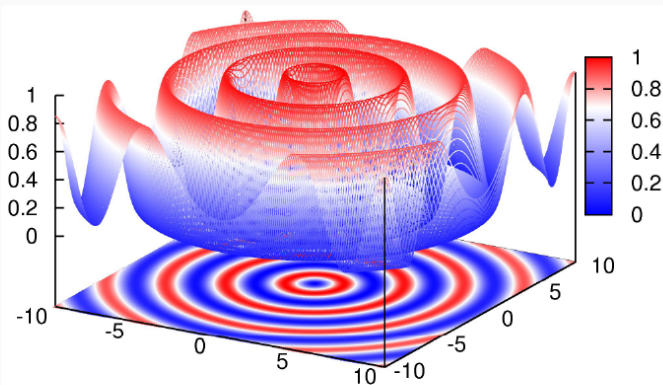
- Simply “repair” infeasible individuals (say when boundary is a box) before update so that they satisfy the constraint
 1. Caveat : Repairs are dangerous
 - Distribution affected by repair, hurting CMA’s convergence
 2. “Re-repair” mechanisms to prevent divergence are also reported
- Alternatively, penalize the repaired solutions

$$f_{\text{fitness}}(\mathbf{x}) = f(\mathbf{x}_{\text{repaired}}) + \alpha \|\mathbf{x} - \mathbf{x}_{\text{repaired}}\|^2$$

Comparing CMA against GA

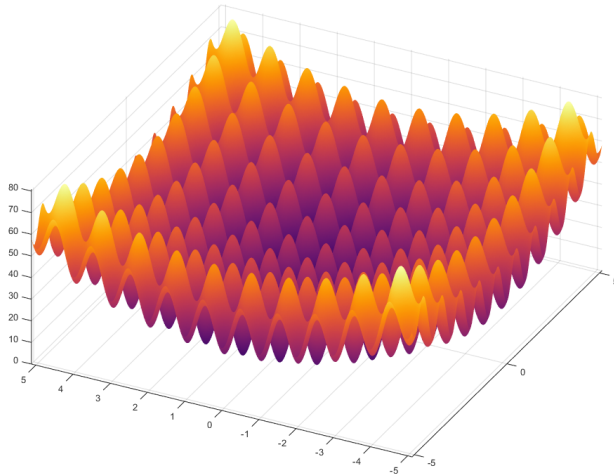
Optimization on smooth functions

- Two dimensional, C^∞ functions $f(\mathbf{x}) : (\mathbb{R}^2, \mathbb{R}, f, \leq)$
- shifted Schaffer function (optima in the middle well)



CMAes vs GA-setup

- shifted Rastrigin function (optima in the middle well)



Comparison between functions³

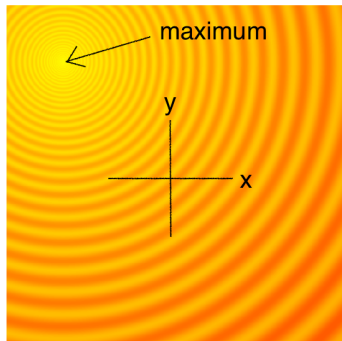


Figure 2: Schaffer-setup

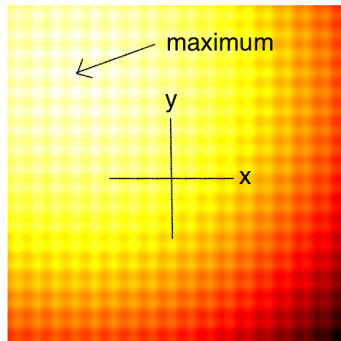


Figure 3: Rastrigin-setup

Lighter region indicates smaller values

³Otoro

Simple ES

Scheme

- **Sampling** : $\mathbf{z}_i \sim \mathcal{N}(\mathbf{m}, \mathbf{C})$
- **Mean-update** : $\mathbf{m} \leftarrow \mathbf{z}_{1:\lambda}$
- **Covariance-update** : $\mathbf{C} = \begin{bmatrix} \sigma_x^2 & \sigma_x \sigma_y \\ \sigma_x \sigma_y & \sigma_y^2 \end{bmatrix}$ σ_x, σ_y are fixed.
- No other updates (on path etc.)

Legend

- **Green** : Tracks the mean \mathbf{m} .
- **Blue** : Tracks the sampled solutions at generation g .
- **Red** : Tracks the best individual so far.

Results

Simple Evolution strategy from Otoro shown for 20 generations

Convergence

- What do you expect for general problems?
-

Rate of convergence

- Is this fast/slow convergence?
-

Number of function evaluations?

- High? Low? Not bad?
-

Convergence

- What do you expect for general problems?
- Will get stuck–lack of diversity, keeps only best population
(See rastrigin, which temporarily gets stuck)
- Heavy parameter dependence too

Rate of convergence

- Is this fast/slow convergence?
- Slow–no history information

Number of function evaluations?

- High? Low? Not bad?
- Decent–but no promises for real life black-box optimization problems

Scheme

- **Environmental selection** : Keep only best 10%
- **Sampling** : Crossover from parents selected above with $p_c = 1$
- **Crossover** : Select two parents, obtain x or y from either parent with 0.5 probability (two coin tosses)
- **Mutation** : Introduce Gaussian noise with fixed σ
- No other updates (on path etc.)

Legend

- **Green** : Tracks the elites from prior generation g .
- **Blue** : Offsprings from candidate solutions.
- **Red** : Tracks the best individual so far.

Convergence

- What do you expect for general problems?
-

Rate of convergence

- Is this fast/slow convergence?
-

Number of function evaluations?

- High? Low? Not bad?
-

Convergence

- What do you expect for general problems?
- Will get stuck—lack of diversity, keeps only elitist population
- Heavy parameter dependence
- Tracks modality well (for both Schaffer and Rastrigin)

Rate of convergence

- Is this fast/slow convergence?
- Slower than simple ES

Number of function evaluations?

- High? Low? Not bad?
- High

Can you spot the updates?

- m update (fairly obvious)
- Step size update
 - Path update
- Covariance matrix update
 - Rank μ updates
 - Rank 1 update (Path update)

Convergence

- What do you expect for general problems?
- **Good** for problems of “moderate” dimensions

Rate of convergence

- Is this fast/slow convergence?
- **Fast** (Approximately brackets minima in $\mathcal{O}(n)$ functional evaluations)

Number of function evaluations?

- High? Low? Not bad?
- Low (same as above)

CMAes-Some interesting videos

- Mario
<https://www.youtube.com/watch?v=0iipyd7Gi70>
- Rastrigin :
<https://www.youtube.com/watch?v=aP31Q7o2UGU>
- Biped:
<https://www.youtube.com/watch?v=l0aWv0A9cb4>
- Robot Invivo:
<https://www.youtube.com/watch?v=trR2Gc1tLzg>
- Robot invitro:
<https://www.youtube.com/watch?v=fjTd06L-9bQ>
- Knifefish :
<https://www.youtube.com/watch?v=3XjgZbs0t2g>
- <https://blog.openai.com/evolution-strategies/>