DERIVATIVE FREE OPTINIZATION 2024/2025.





Pure Random Search (PRS)

 $f: f = 1, i \xrightarrow{n} \longrightarrow \mathbb{R}$ $(x = b_{1}, \dots, x_{n}) \longrightarrow f(x)$ ASSume <u>PRS</u>: Initialize xbest = Unif ([-1,13^h) (uniforma) WHILE NOT HAPPY (while stop criterion not met) Sample X~ Unif (C-1,13ⁿ) If $f(x) \leq f(xbest)$ Xbest <- X

Does this algorithm converge : Yes under mild assumptions on f (need to have "volume" in a reighborhood of global prime





f(x)= 11×1100 = max (kul, -.; kul)





Vi~ Unif([-1])

Xt: best solution at iteration t

 $f(Xt) = \min \left\{ f(u_{i}), \dots, f(u_{t}) \right\}$



By induction.

 $\lim_{t \to +\infty} \frac{\mathbb{P}(1|X + 1|_{00} \ge \varepsilon)}{1 + 1} = 0$ Pour texo

La give ev in probability.













Loore otherwite.

TE: time it takes to win this game.

Given a game with 2 outcomes win with pobe p and look with

proba (1-p), where the outcome is sampled randomly and independently

[example : flip a coin], the time it takes to win a game is

distributed according to a geometric distribution.

 $E[T_2] = \frac{1}{p}$

Back to PRS. $p=\operatorname{IP}("win")=\operatorname{IP}(||V+|| \leq \varepsilon)=\varepsilon^m$



$\mathbb{E}(T_{\mathcal{E}}) = \frac{1}{\varepsilon^{m}}$

Is it fast? No



The algorithm is "blind", does not take into account the

information gæhered on f, through the sampling of points to sample "better" solutions.



Part II: Algorithms

Deterministic Algorithms

Quasi-Newton with estimation of gradient (BFGS) [Broyden et al. 1970] Simplex downhill [Nelder and Mead 1965]

Pattern search, Direct Search [Hooke and Jeeves 1961]

Trust-region/Model Based methods (NEWUOA, BOBYQA) [Powell, 06,09]

Stochastic (randomized) search methods

Evolutionary Algorithms (continuous domain)
Differential Evolution [Storn, Price 1997]
Particle Swarm Optimization [Kennedy and Eberhart 1995]
Evolution Strategies, CMA-ES [Rechenberg 1965, Hansen, Ostermeier 2001]
Estimation of Distribution Algorithms (EDAs) [Larrañaga, Lozano, 2002]
Cross Entropy Method (same as EDAs) [Rubinstein, Kroese, 2004]
Genetic Algorithms [Holland 1975, Goldberg 1989]

Simulated Annealing [Kirkpatrick et al. 1983]

A Generic Template for Stochastic Search

Define $\{P_{\theta} : \theta \in \Theta\}$, a family of probability distributions on \mathbb{R}^{n}

Generic template to optimize $f : \mathbb{R}^n \to \mathbb{R}$

Initialize distribution parameter θ , set population size $\lambda \in \mathbb{N}$

- While not terminate 1. Sample $x_1, ..., x_{\lambda}$ according to P_{θ}
 - 2. Evaluate x_1, \ldots, x_{λ} on f
 - 3. Update parameters $\theta \leftarrow F(\theta, x_1, \dots, x_{\lambda}, f(x_1), \dots, f(x_{\lambda}))$

the update of θ should drive P_{θ} to concentrate on the optima of f

To obtain an optimization algorithm we need:

• to define $\{P_{\theta}, \theta \in \Theta\}$ • to define F the update function of θ

Which probability distribution to sample candidate solutions?

Normal distribution - 1D case



General case

• Normal distribution $\mathcal{N}(\boldsymbol{m}, \sigma^2)$

probability density of the 1-D standard normal distribution $\mathcal{N}(0,1)$

(expected (mean) value, variance) = (0,1)

$$p(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$

(expected value, variance) = $(\boldsymbol{m}, \sigma^2)$ density: $p_{\boldsymbol{m},\sigma}(x) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(x-\boldsymbol{m})^2}{2\sigma^2}\right)$

- A normal distribution is entirely determined by its mean value and variance
- The family of normal distributions is closed under linear transformations: if X is normally distributed then a linear transformation aX + b is also normally distributed

• Exercice: Show that
$$m + \sigma \mathcal{N}(0, 1) = \mathcal{N}(m, \sigma^2)$$

m + o Wb,1) is normally distributed

We only need to identify its mean and variance:

 $E\left(m + \sigma W(o, 1)\right) = m + \sigma E(W_{0, 1}) - m$ by linenity =0 or E Var(m + \sigma W(o, 1)) $= \mathbb{E}\left(\left(mr\sigma \mathcal{N}(q_1) - m\right)^2\right) = \mathbb{E}\left(\sigma^2 \mathcal{N}(0, \Lambda)\right)$ $= \sigma^2 \mathbb{E}(\mathbb{W}(0,1)^2) = \sigma^2$ = 1

= m $\sigma \sigma \mathcal{N}(G,1) \simeq \mathcal{N}(m, \sigma^2)$





Generalization to n Variables: Independent Case

Assume X1 ~
$$\mathcal{N}(\mu_1, \sigma_1^2)$$
 denote its density $p(x_1) = \frac{1}{Z_1} \exp\left(-\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2\right)$
Assume X2~ $\mathcal{N}(\mu_2, \sigma_2^2)$ denote its density $p(x_2) = \frac{1}{Z_2} \exp\left(-\frac{1}{2\sigma_2^2}(x_2 - \mu_2)^2\right)$

Assume X1 and X2 are **independent**, then (X1,X2) is a Gaussian vector with

$$p(x_1, x_2) =$$

Generalization to n Variables: Independent Case

Assume X1 ~
$$\mathcal{N}(\mu_1, \sigma_1^2)$$
 denote its density $p(x_1) = \frac{1}{Z_1} \exp\left(-\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2\right)$
Assume X2~ $\mathcal{N}(\mu_2, \sigma_2^2)$ denote its density $p(x_2) = \frac{1}{Z_1} \exp\left(-\frac{1}{2\sigma_2^2}(x_2 - \mu_2)^2\right)$

Assume X1 and X2 are **independent**, then (X1,X2) is a Gaussian vector with

$$p(x_1, x_2) = p(x_1)p(x_2) = \frac{1}{Z_1 Z_2} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

with $x = (x_1, x_2)^T$ $\mu = (\mu_1, \mu_2)^T$ $\Sigma = \begin{pmatrix} \sigma_1^2 & 0\\ 0 & \sigma_2^2 \end{pmatrix}$

Generalization to n Variables: Independent Case

Assume X1 ~
$$\mathcal{N}(\mu_1, \sigma_1^2)$$
 denote its density $p(x_1) = \frac{1}{Z_1} \exp\left(-\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2\right)$
Assume X2~ $\mathcal{N}(\mu_2, \sigma_2^2)$ denote its density $p(x_2) = \frac{1}{Z_1} \exp\left(-\frac{1}{2\sigma_2^2}(x_2 - \mu_2)^2\right)$

Assume X1 and X2 are **independent**, then (X1,X2) is a Gaussian vector with $(-\frac{1}{2}) \left(\frac{x_n - y_n}{x_n - y_n}\right)^2 + \frac{x_n - y_n}{x_n - y_n}$

$$p(x_1, x_2) = p(x_1)p(x_2) = \frac{1}{Z_1 Z_2} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

with $x = (x_1, x_2)^T$ $\mu = (\mu_1, \mu_2)^T$ $\Sigma = \begin{pmatrix} \sigma_1^2 & 0\\ 0 & \sigma_2^2 \end{pmatrix}$



Gaussian Vector - Multivariate Normal Distribution

A random vector $X = (X_1, ..., X_n) \in \mathbb{R}^n$ is a Gaussian vector (or multivariate normal) if and only if for all real numbers $a_1, ..., a_n$, the random variable $a_1X_1 + ... + a_nX_n$ has a normal distribution.

Gaussian Vector - Multivariate Normal Distribution

A random variable following a 1-D normal distribution is determined by its mean value m and variance σ^2 .

In the *n*-dimensional case it is determined by its mean vector and covariance matrix

Covariance Matrix

If the entries in a vector $\mathbf{X} = (X_1, \dots, X_n)^T$ are random variables, each with finite variance, then the covariance matrix Σ is the matrix whose (i, j) entries are the covariance of (X_i, X_j)

$$\Sigma_{ij} = \operatorname{cov}(X_i, X_j) = \operatorname{E}\left[(X_i - \mu_i)(X_j - \mu_j)\right]$$

where $\mu_i = E(X_i)$. Considering the expectation of a matrix as the expectation of each entry, we have $\sum_{i=1}^{n} \sum_{j=1}^{n} \sqrt{\alpha_i(X_i)}$

$$\boldsymbol{\Sigma} = \mathrm{E}[(\boldsymbol{X} - \boldsymbol{\mu})(\boldsymbol{X} - \boldsymbol{\mu})^{\mathsf{T}}]$$

 Σ is symmetric, positive definite

Density of a n-dimensional Gaussian vector $\mathcal{N}(m, C)$:

$$p_{\mathcal{N}(m,C)}(x) = \frac{1}{(2\pi)^{n/2} |C|^{1/2}} \exp\left(-\frac{1}{2}(x-m)^{\mathsf{T}}C^{-1}(x-m)\right)$$

(cl=det(c)

The mean vector *m*:

determines the displacement

is the value with the largest density

the distribution is symmetric around the mean

$$\mathcal{N}(m, C) = m + \mathcal{N}(0, C)$$

The covariance matrix:

determines the geometrical shape (see next slides)



Consider a Gaussian vector $\mathcal{N}(m, C)$, remind that lines of equal densities are given by:

$$\{x \mid \Delta^2 = (x - m)^T C^{-1} (x - m) = \text{cst}\}\$$

Decompose $C = U \Lambda U^{T}$ with U orthogonal, i.e.

$$C = \begin{pmatrix} u_1 & u_2 \\ | & | \end{pmatrix} \begin{pmatrix} \sigma_1^2 & 0 \\ 0 | & \sigma_2^2 \end{pmatrix} \begin{pmatrix} u_1 & - \\ u_2 & - \end{pmatrix}$$

Let $Y = U^{\top}(x - m)$, then in the coordinate system, (u1,u2), the lines of equal densities are given by





...any covariance matrix can be uniquely identified with the iso-density ellipsoid $\{x \in \mathbb{R}^n | (x - m)^T C^{-1} (x - m) = 1\}$



where I is the identity matrix (isotropic case) and D is a diagonal matrix (reasonable for separable problems) and $\mathbf{A} \times \mathcal{N}(\mathbf{0}, \mathbf{I}) \sim \mathcal{N}(\mathbf{0}, \mathbf{A}\mathbf{A}^{\mathrm{T}})$ holds for all A.

Evolution Strategies

New search points are sampled normally distributed

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

as perturbations of *m*,

where
$$oldsymbol{x}_i,oldsymbol{m}\in\mathbb{R}^n,\ \sigma\in\mathbb{R}_+,\ oldsymbol{C}\in\mathbb{R}^{n imes n}$$

 $oldsymbol{C}\in\mathbb{R}^{n imes n}$
 $oldsymbol{X}_i\sim\mathbb{W}(oldsymbol{m},\ \sigma^2oldsymbol{C}$



where

- the mean vector $\boldsymbol{m} \in \mathbb{R}^n$ represents the favorite solution
- the so-called step-size $\sigma \in \mathbb{R}_+$ controls the step length
- ► the covariance matrix C ∈ ℝ^{n×n} determines the shape of the distribution ellipsoid

here, all new points are sampled with the same parameters

Evolution Strategies

New search points are sampled normally distributed

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

- In fact, the covariance matrix of the sampling distribution is $\sigma^2 \mathbb{C}$ but it is convenient to refer to \mathbb{C} as the covariance matrix (it is a covariance matrix but not of the sampling distribution)
 - the mean vector $m \in \mathbb{R}^n$ represents the favorite solution
 - the so-called step-size $\sigma \in \mathbb{R}_+$ controls the step length
 - ► the covariance matrix C ∈ ℝ^{n×n} determines the shape of the distribution ellipsoid

here, all new points are sampled with the same parameters

How to update the different parameters m, σ, \mathbf{C} ?

1. Adapting the mean *m*

- 2. Adapting the step-size σ
- **3.** Adapting the covariance matrix *C*

Update the Mean: a Simple Algorithm the (1+1)-ES

Notation and Terminology:

one solution kept from one iteration to the next



one new solution (offspring) sampled at each iteration

The + means that we keep the best between current solution and new solution, we talk about *elitist* selection

(1+1)-ES algorithm (update of the mean)

sample one candidate solution from the mean ${\boldsymbol{m}}$

 $\mathbf{x} = \mathbf{m} + \sigma \mathcal{N}(0, \mathbf{C})$

if **x** is better than **m** (i.e. if $f(\mathbf{x}) \leq f(\mathbf{m})$), select **m**

 $\mathbf{m} \leftarrow \mathbf{x}$