

stochastics and probability

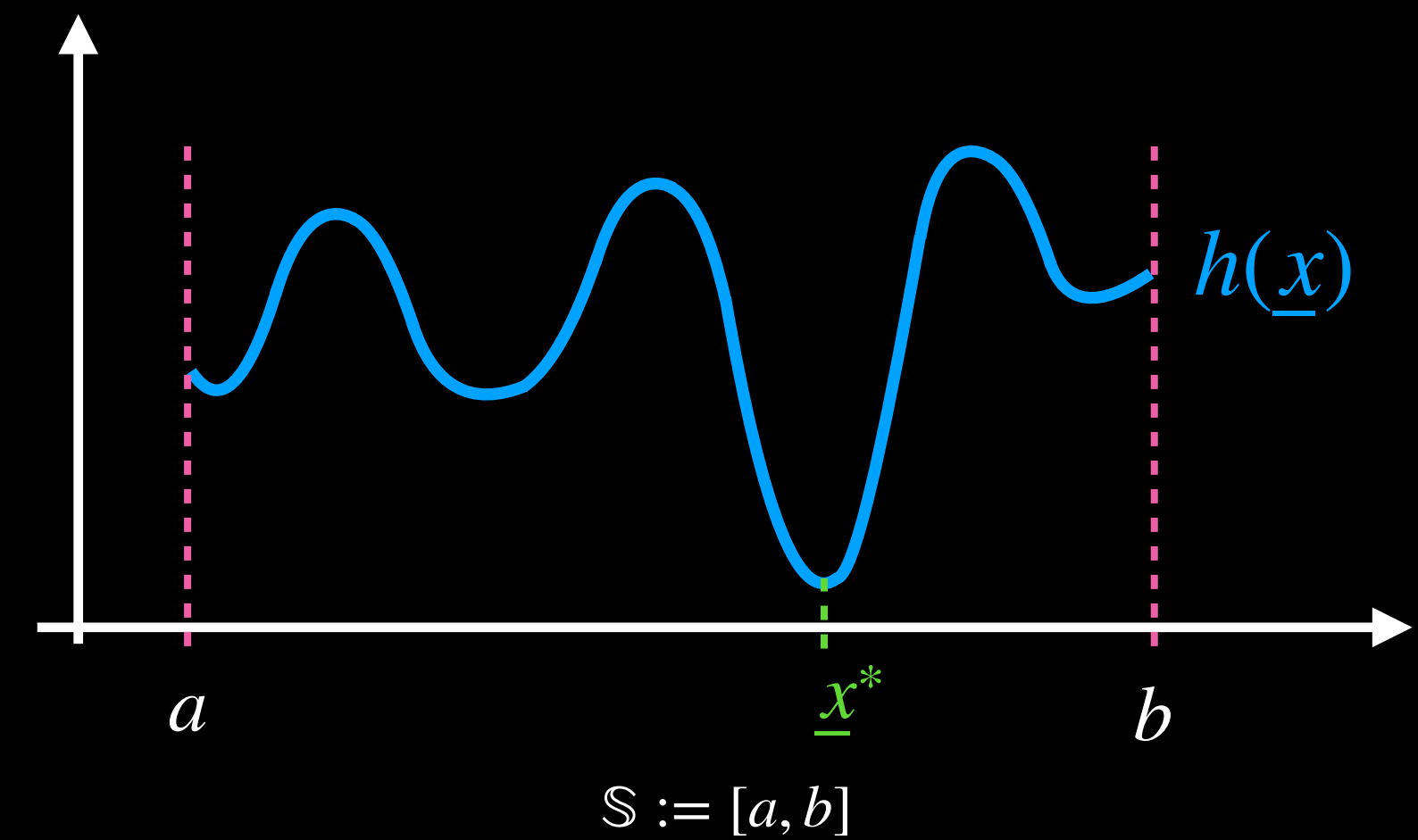
Lecture 10

Dr. Johannes Pahlke

optimization

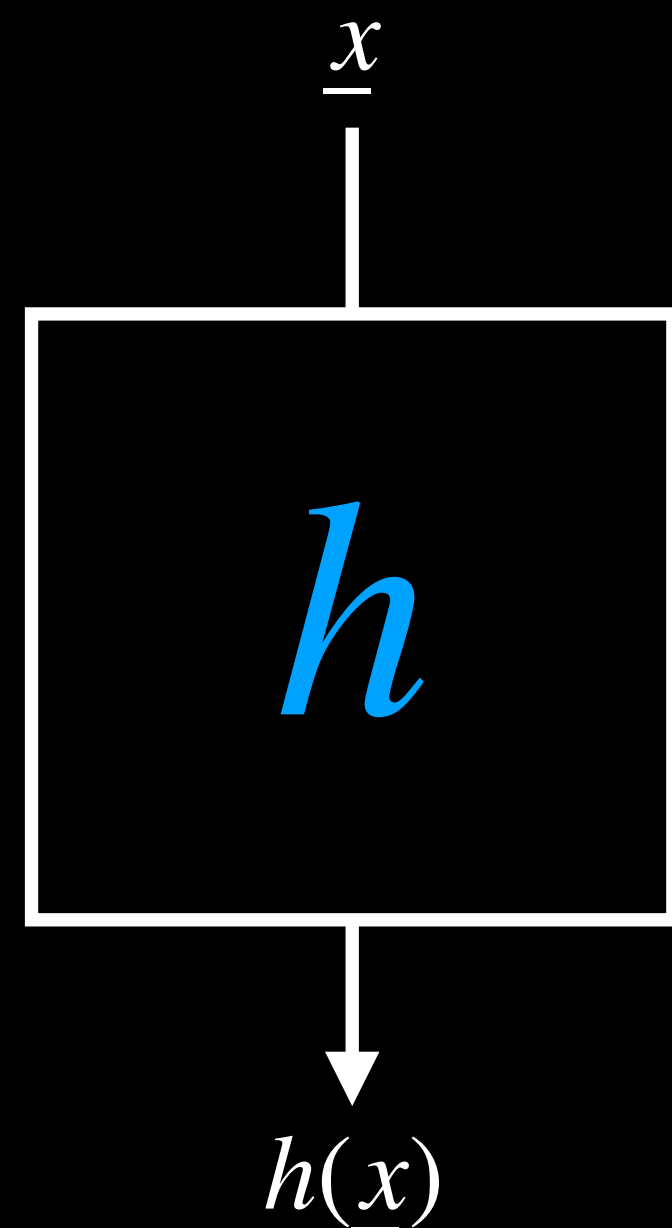
Goal:

$$\underline{x}^* := \arg \min_{\underline{x} \in \mathbb{S}} (h(\underline{x})) \quad \left(\longleftrightarrow \quad h(\underline{x}^*) = \min_{\underline{x} \in \mathbb{S}} (h(\underline{x})) \right)$$



Challenges:

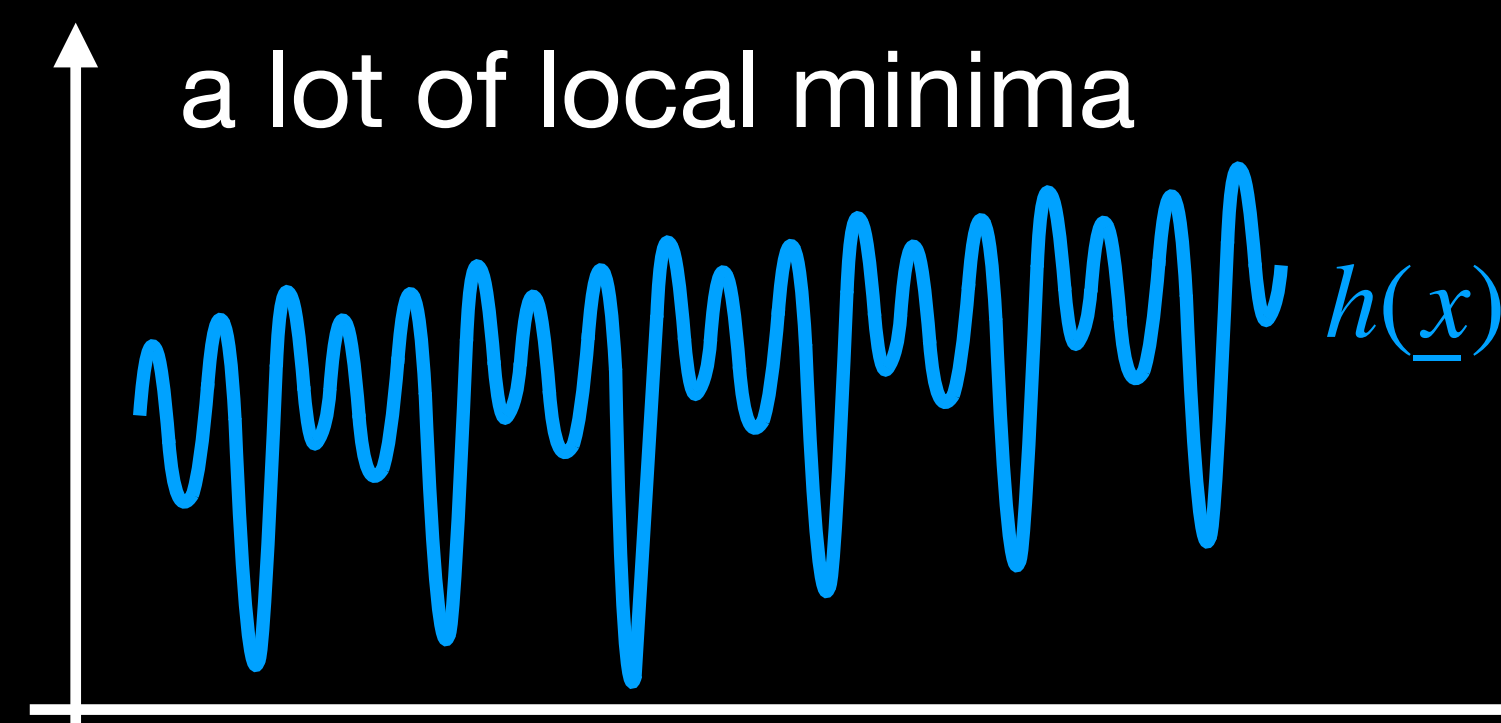
black box



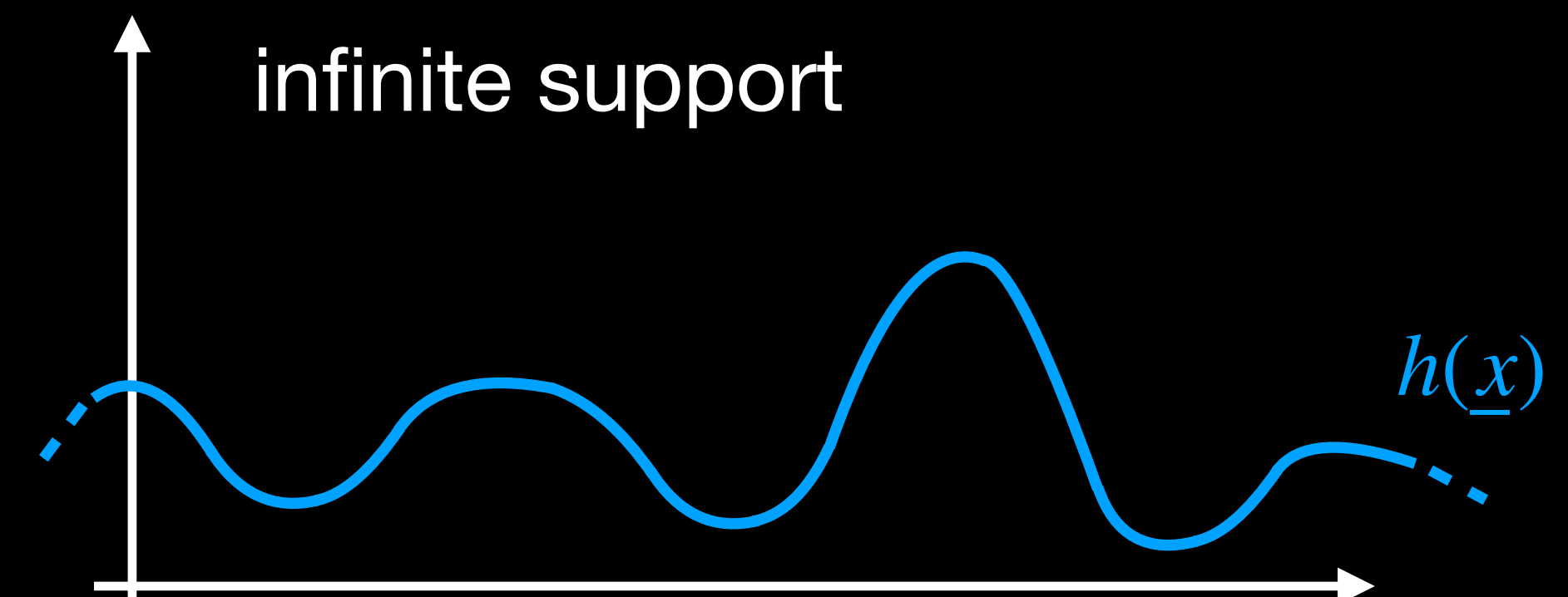
high dimensional

$$h(x_1, x_2, x_3, \dots, x_{100})$$

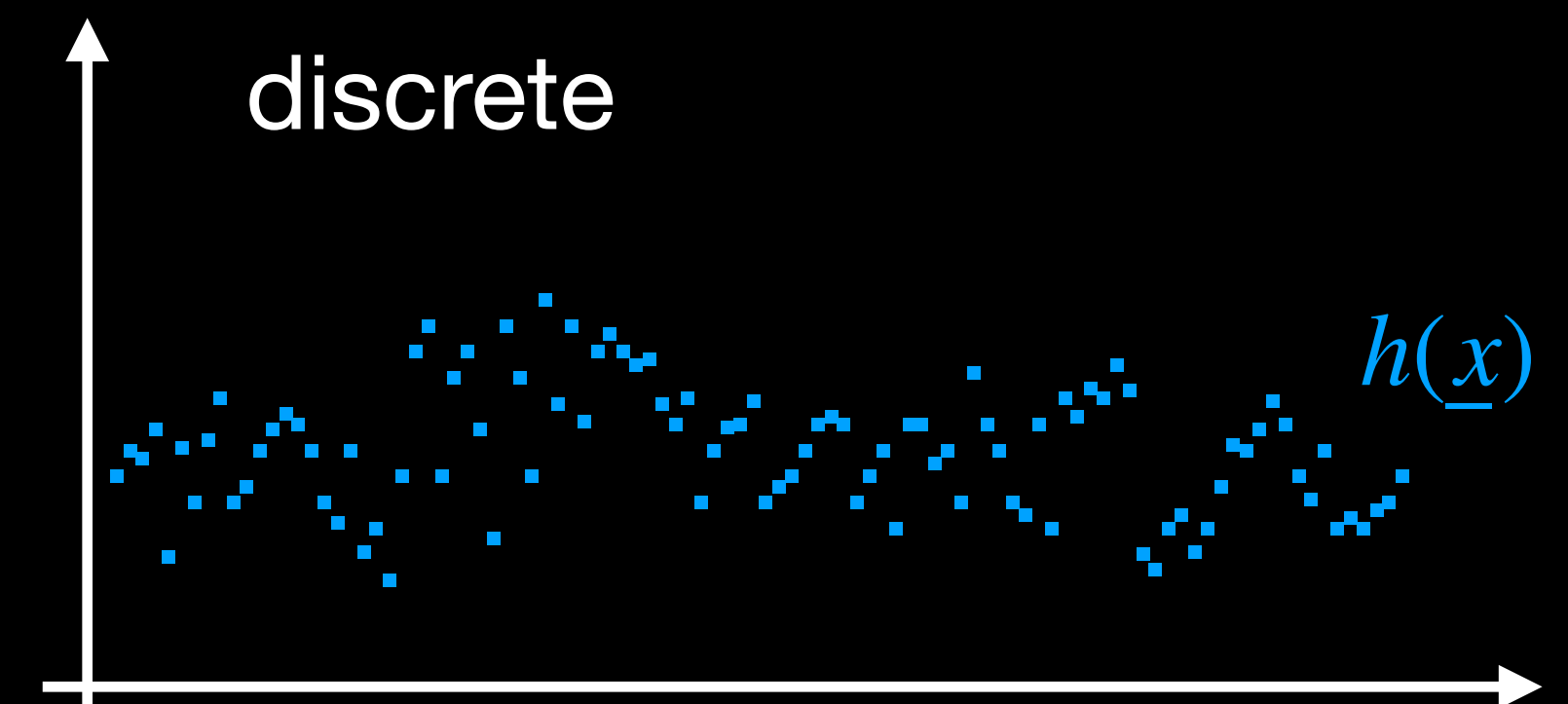
function evaluation are expensive



infinite support



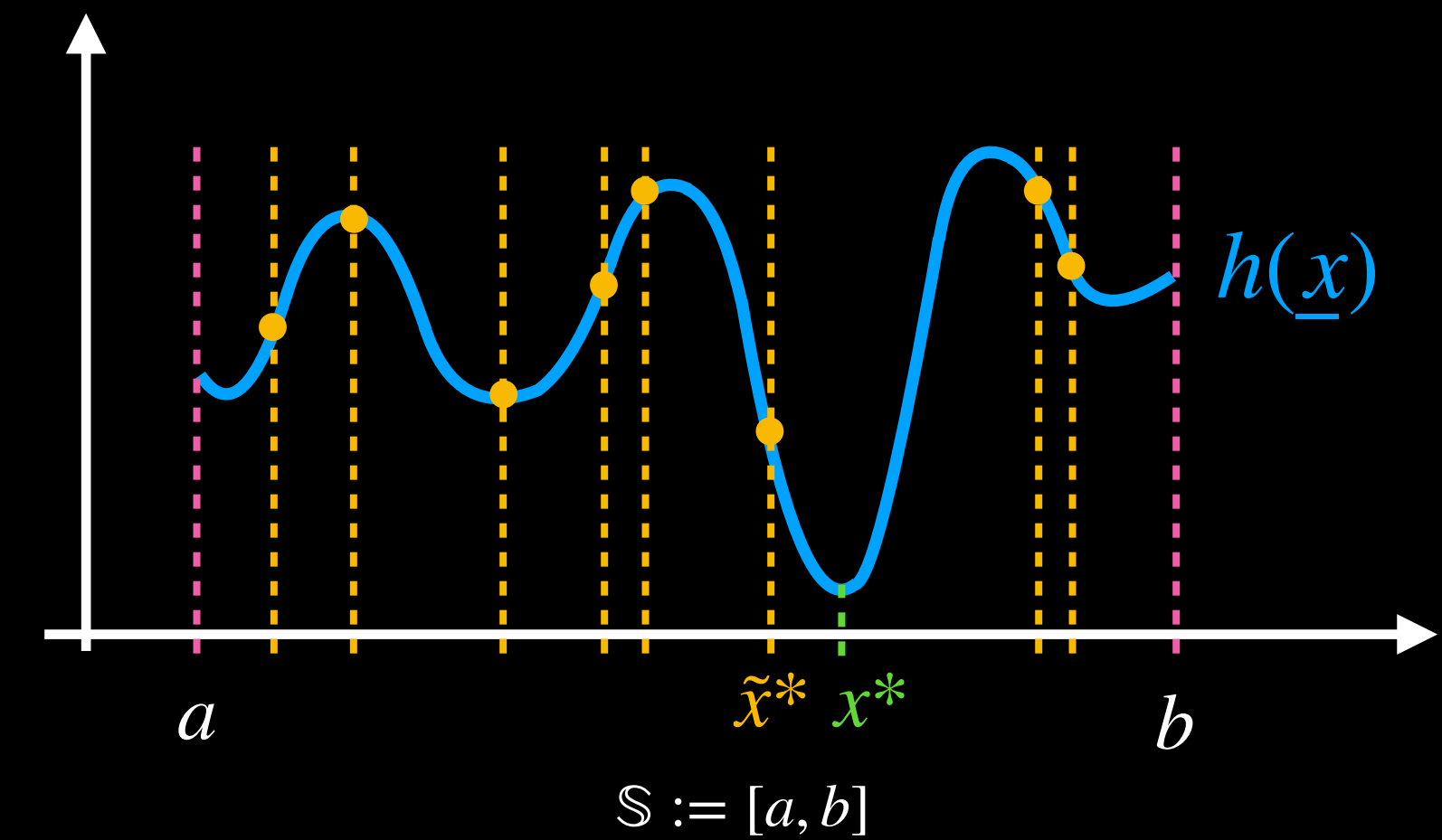
discrete



algorithm 1: stochastic exploration

$$\underline{x}^* := \arg \min_{\underline{x} \in \mathbb{S}} (h(\underline{x})) \approx \arg \min (h(\underline{X}_1), \dots, h(\underline{X}_n)) =: \tilde{\underline{X}}_n^* \quad \underline{X}_t : \Omega \rightarrow \mathbb{S} \subseteq \mathbb{R}^d$$

$$\underline{X}_t \sim \mathcal{U}(\mathbb{S})$$



global convergence:

$$\underline{x}^* = \lim_{n \rightarrow \infty} \tilde{\underline{X}}_n^*$$

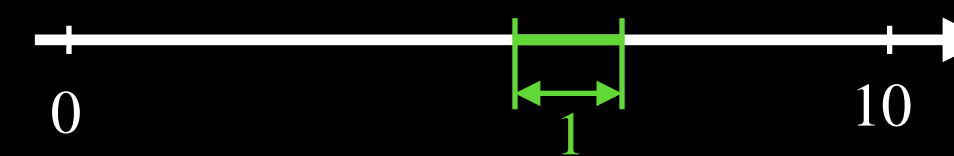
properties:

- very simple
- converges to \underline{x}^*

useful if:

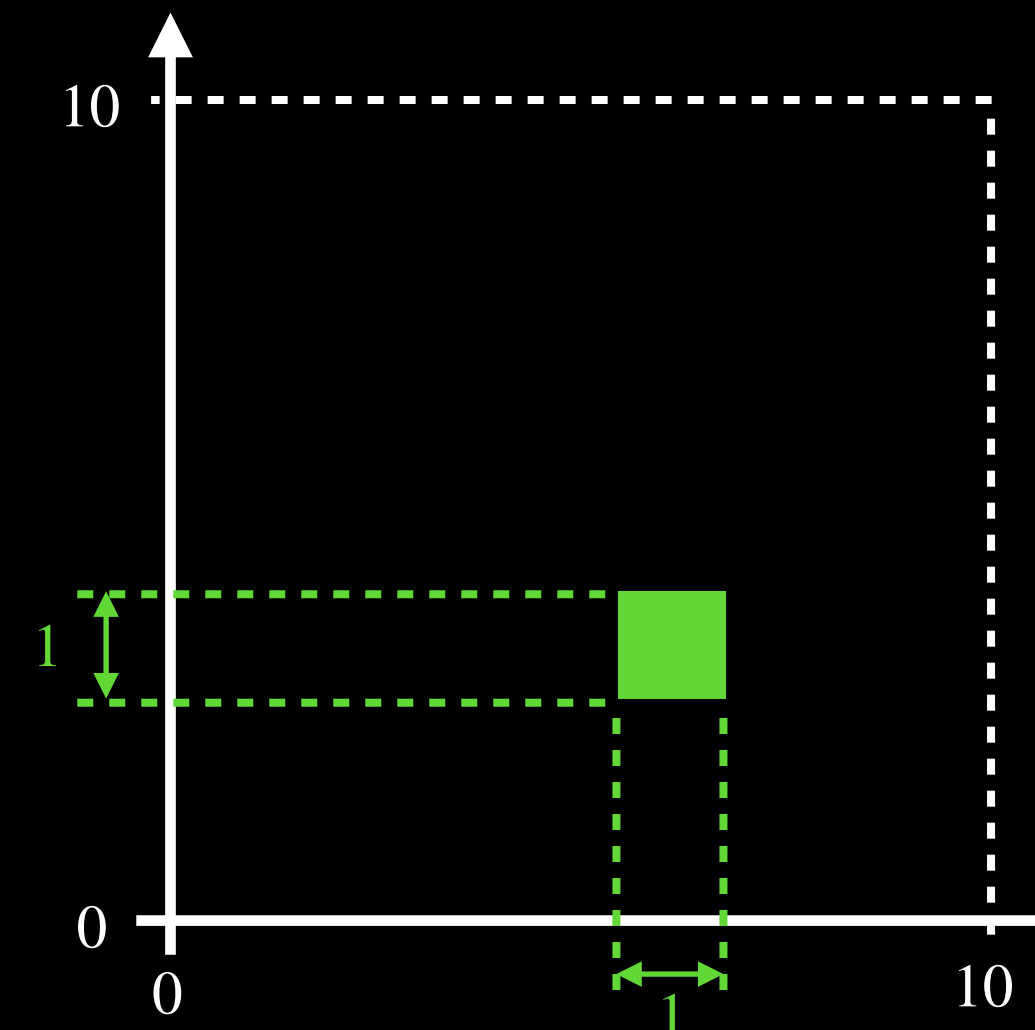
- \mathbb{S} is low-dimensional and bounded
- (h is cheap to evaluate for higher dimensions)

1d:



$$P(X \in \text{—}) = \frac{1}{10}$$

2d:



$$P(X \in \blacksquare) = \frac{1}{10^2}$$

$$\text{error} \propto n^{-\frac{1}{d}}$$

algorithm 2: stochastic descent

$$\underline{x}^* := \arg \min_{\underline{x} \in \mathbb{S}} h(\underline{x})$$

gradient descent: $\underline{x}_{t+1} = \underline{x}_t - \alpha_t \nabla h(\underline{x}_t)$

$t \in \{1, \dots, n\}$ steps
 $\alpha_t > 0$ step size
 $-\nabla h(\underline{x}_t)$ vector of steepest decrease
 \underline{x}_0 starting point

gradient:

$$\nabla h(\underline{x}) := \begin{pmatrix} \frac{\partial}{\partial x_1} h(\underline{x}) \\ \vdots \\ \frac{\partial}{\partial x_d} h(\underline{x}) \end{pmatrix}, \quad \frac{\partial}{\partial x_1} h(\underline{x}) \approx \frac{h \begin{pmatrix} x_1 + \Delta x \\ x_2 \\ \vdots \\ x_d \end{pmatrix} - h \begin{pmatrix} x_1 - \Delta x \\ x_2 \\ \vdots \\ x_d \end{pmatrix}}{2\Delta x}$$

→ approximating $\nabla h(\underline{x})$ has $2d$ evaluations of h (expensive)

stochastic descent:

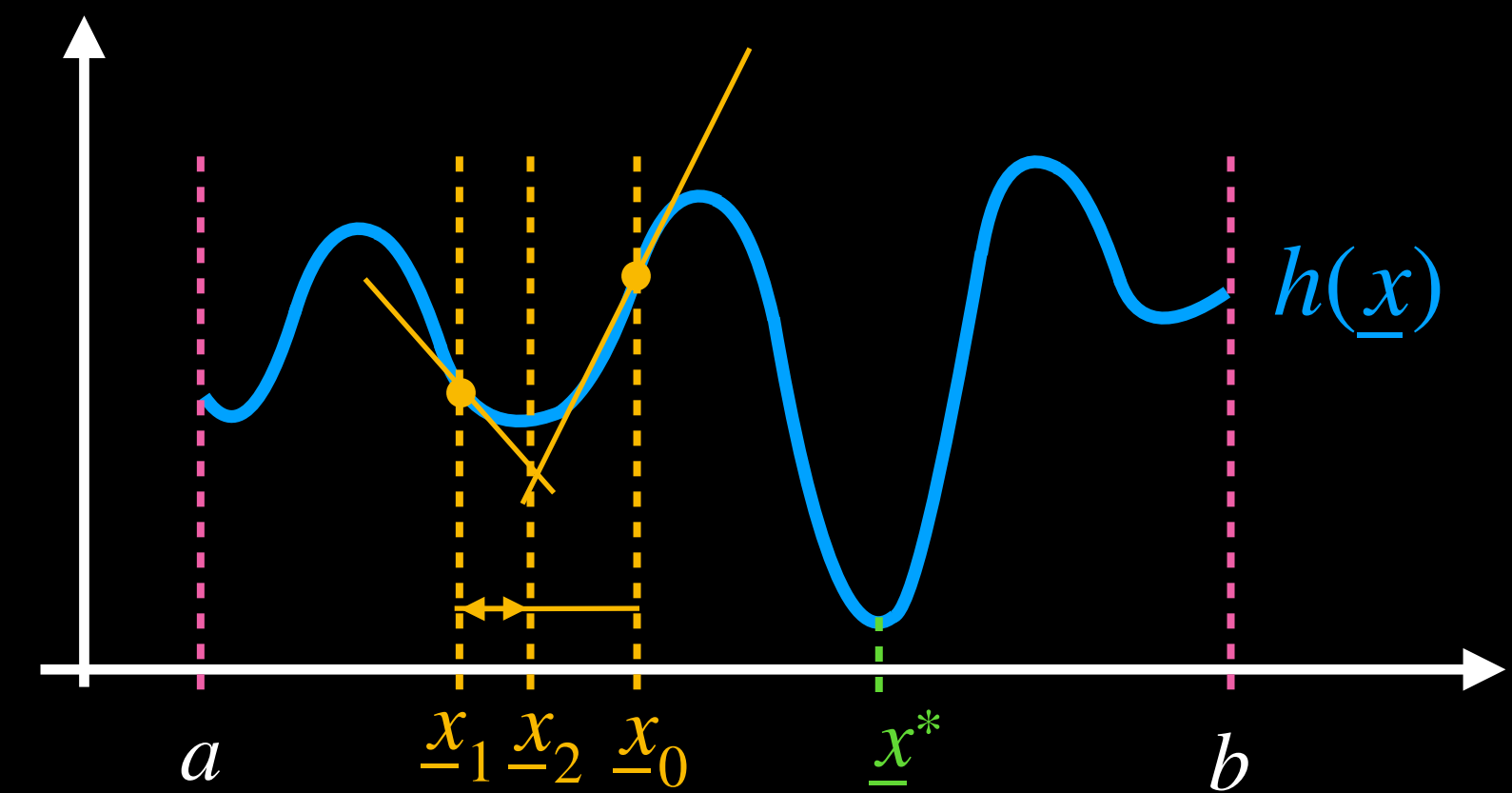
$$\underline{X}_{t+1} = \underline{X}_t + \frac{\alpha_t}{2\beta_t} \left(\nabla h(\underline{X}_t) U_t - \nabla h(\underline{X}_{t-1}) U_{t-1} \right)$$

$$\underline{U}_t \sim \mathcal{U}(\mathbb{S})$$

\mathbb{S} be the d -dimensional unit sphere
 β_t sampling radius

properties:

- simple
- no global convergence guarantee
- local convergence if $\lim_{n \rightarrow \infty} \alpha_n = 0$ and $\lim_{n \rightarrow \infty} \frac{\alpha_n}{\beta_n} = \text{const.}$
- converges fast ($\propto \frac{1}{n}$)
- needs uniform random numbers on the unit sphere
- need to choose $\alpha_t, \beta_t, \underline{X}_0$



dot product: $\underline{x} \cdot \underline{y} = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} \cdot \begin{pmatrix} y_1 \\ \vdots \\ y_d \end{pmatrix} = x_1 y_1 + \dots + x_d y_d$

directional derivative:

$$\nabla h(\underline{x}) \cdot \frac{\underline{y}}{|\underline{y}|} \approx \frac{\Delta h(\underline{x}, \underline{y})}{2|\underline{y}|}, \quad \Delta h(\underline{x}, \underline{y}) := h(\underline{x} + \underline{y}) - h(\underline{x} - \underline{y})$$

→ approximation has 2 evaluations of h (cheaper)

useful if:

- \mathbb{S} is high-dimensional and unbounded
- h is convex or \underline{X}_0 is close to \underline{x}^*

algorithm 3: random pursuit

$$\underline{x}^* := \arg \min_{\underline{x} \in \mathbb{S}} (h(\underline{x}))$$

$$\underline{X}_{t+1} = \underline{X}_t + \text{linesearch}(\underline{X}_t, \underline{U}_t) \quad \underline{U}_t \sim \mathcal{U}(\mathbb{S})$$

\mathbb{S} be the d -dimensional unit sphere

$$\text{linesearch}_h(\underline{X}_t, \underline{U}_t) := \arg \min_{\beta} \left(h(\underline{X}_t + \beta \cdot \underline{U}_t) \right) \cdot \underline{U}_t \quad \beta \text{ is distance between } \underline{X}_t \text{ and } \underline{X}_{t+1}$$

(finds the vector to the point with minimal value on the line that goes through \underline{X}_t in the direction \underline{U}_t)

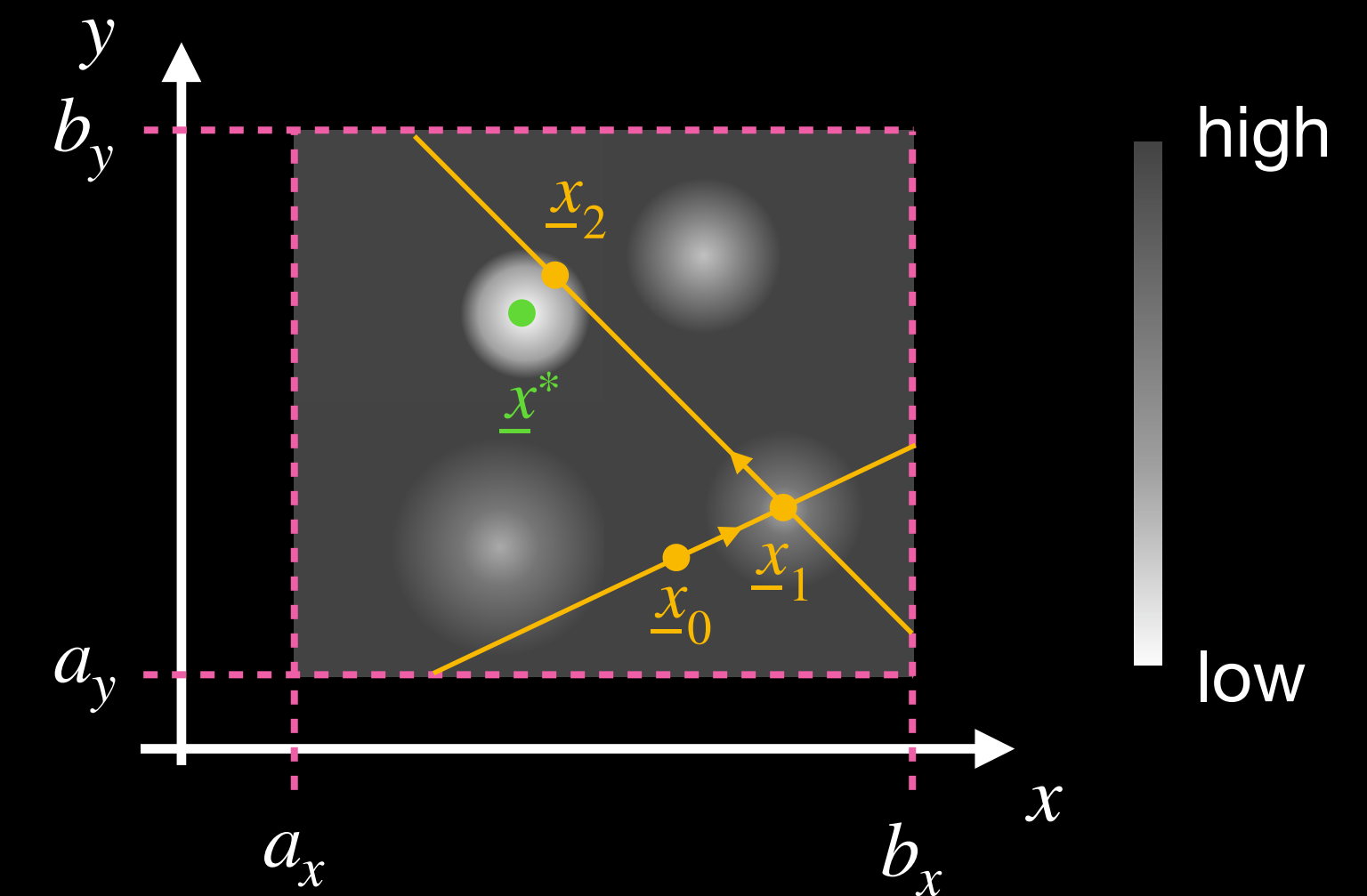
→ reduction to 1d-optimisation problem

properties:

- global convergence guarantee
- converges fast
- needs uniform random numbers (on the unit sphere)

useful if:

- \mathbb{S} is high-dimensional and bounded



algorithm 4: simulated annealing

$$\underline{x}^* := \arg \min_{\underline{x} \in \mathbb{S}} (h(\underline{x}))$$

$$\underline{X}_{t+1} := \begin{cases} \underline{X}_t + \underline{U}_t^\beta & \text{if } U_t < e^{-\frac{h(\underline{X}_t) - h(\underline{X}_t + \underline{U}_t^\beta)}{T_t}} \\ \underline{X}_t & \text{else} \end{cases} \quad \begin{aligned} U_t &\sim \mathcal{U}(0,1) \\ \underline{U}_t^\beta &\sim \mathcal{U}(-\frac{1}{2}\underline{\beta}, \frac{1}{2}\underline{\beta}) \end{aligned}$$

$T_t > 0$ is the "temperature" (parameter for the chance to accept worse points)

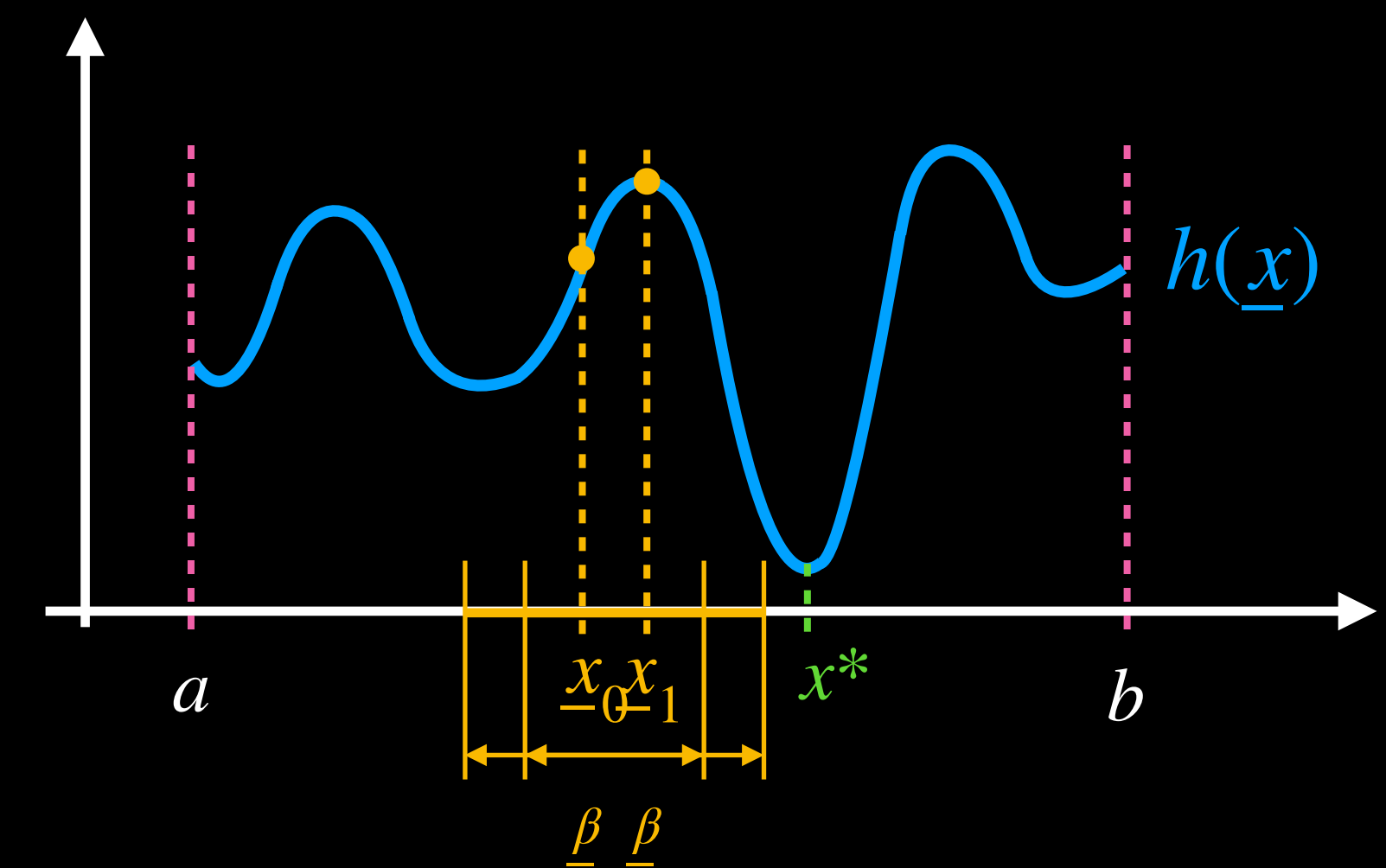
(example for a cooling scheme: $T_{t+1} := 0.95 T_t$)

properties:

- simple
- convergence to a local minimum
- global convergence for the "right" cooling scheme
- can go out of local minima again

useful if:

- \mathbb{S} is high-dimensional
- finding the global optimum is more important than finding the exact position of a (local/ global) optimum
- h has a lot of local minima
- h is discrete



algorithm 5: evolutionary algorithms

$$\underline{x}^* := \arg \min_{\underline{x} \in \mathbb{S}} (h(\underline{x}))$$

(1+1)-ES: (ES = Evolution Strategy)

$$\underline{Y}_t = \underline{X}_t + \underline{N}_t \quad \underline{N}_t \sim \mathcal{N}(\underline{0}, \underline{\Sigma})$$

$$\underline{X}_{t+1} = \arg \min (h(\underline{Y}_t), h(\underline{X}_t))$$

\underline{Y}_t child
 \underline{X}_t parent
 \underline{N}_t mutation
 $\mathcal{N}(\underline{0}, \underline{\Sigma})$ multivariate Gaussian
 $\underline{\Sigma}$ co-variance matrix

properties:

- global convergence possible
- (can go out of local minima)

(1+ λ)-ES:

$$\underline{Y}_{t,k} = \underline{X}_t + \underline{N}_{t,k} \quad \underline{N}_{t,k} \sim \mathcal{N}(\underline{0}, \underline{\Sigma})$$

$$\underline{X}_{t+1} = \arg \min (h(\underline{Y}_{t,1}), \dots, h(\underline{Y}_{t,\lambda}), h(\underline{X}_t))$$

useful if:

- \mathbb{S} is high-dimensional
- finding the global optimum is more important than finding the exact position of a (local/ global) optimum
- h has a lot of local minima
- h is discrete

(μ, λ)-ES & ($\mu + \lambda$)-ES :

$$(\underline{Y}'_{t,1}, \dots, \underline{Y}'_{t,\lambda}) = \text{recombination}(\underline{X}_{t,1}, \dots, \underline{X}_{t,\mu})$$

recombination

$$\underline{Y}_{t,k} = \underline{Y}'_{t,k} + \underline{N}_{t,k} \quad \underline{N}_{t,k} \sim \mathcal{N}(\underline{0}, \underline{\Sigma})$$

mutation

$$(\underline{X}_{t+1,1}, \dots, \underline{X}_{t+1,\mu}) = \text{selection}_{(\mu, \lambda)\text{-ES}}(\underline{Y}_{t,1}, \dots, \underline{Y}_{t,\lambda})$$

selection

$$(\underline{X}_{t+1,1}, \dots, \underline{X}_{t+1,\mu}) = \text{selection}_{(\mu + \lambda)\text{-ES}}(\underline{Y}_{t,1}, \dots, \underline{Y}_{t,\lambda}, \underline{X}_{t,1}, \dots, \underline{X}_{t,\mu})$$

improvement by
covariance matrix adaptation (CMA)

$$\underline{\Sigma}_t$$

- scaling per dimension (different diagonal entries)
- correlations (off-diagonal entries)

end