

## ESs II Theoretical Aspects

Algorithms won their first reputation in practice;  
i.e. they did their job  
in finding good solutions  
for difficult tasks

What is the aim of theory, then?

Insight: why and when do algorithms work  
effectivity, robustness

efficiency, convergence velocity (order)

Insight II: how to improve algorithms

Important questions:

$n$ -dependency of necessary effort,  
speedup for parallel implementations

In the following: restriction to real-valued  
parameter optimization

$(1+1)$ ,  $(1+\lambda)$ , and  $(\mu+\lambda)$  ESs

results from Rechenberg, Rappaport, Bäck,  
Rudolph, Beyer, ...

## 7.1 Some Def. from Prob. Th. & Statistics

### 7.1.1 Random Variables, Distribution & Density Functions

Let  $(\Omega, \mathcal{A}, p)$  a probability space

$X: \Omega \rightarrow \mathbb{R}$  a continuous random variable

$\Omega$  a set of elementary events

$\mathcal{A}$  an algebra of the events

$p$  a probability measure,

then  $F_X: \mathbb{R} \rightarrow [0, 1]$  distribution function

$$x \mapsto F_X(x) = p(X \leq x)$$

$f_X: \mathbb{R} \rightarrow \mathbb{R}$  density function

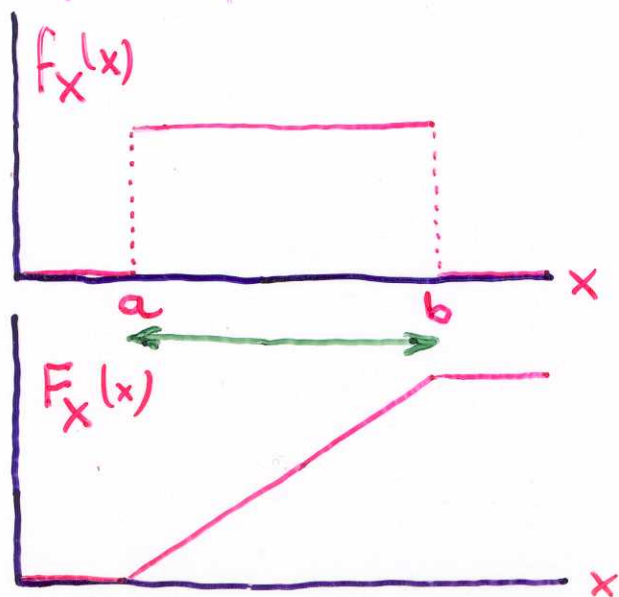
$$F_X(x) = \int_{-\infty}^x f_X(x) dx \quad \forall x \in \mathbb{R}$$

(implicit definition)

or  $f_X(x) = \frac{d}{dx} F_X(x)$

The set  $\{x \in \mathbb{R} \mid f_X(x) > 0\}$  is called support (X)

e.g. uniform random number



$$f_X(x) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b \\ 0 & \text{else} \end{cases}$$

$a < x < b$  support (X)

$$F_X(x) = \begin{cases} 0 & \text{if } x \leq a \\ \frac{x-a}{b-a} & \text{if } a \leq x \leq b \\ 1 & \text{if } x \geq b \end{cases}$$

## 7.1.2 Characteristic Values of p.d.s

### One dimensional case

$$\xi = E[X] := \int_{-\infty}^{+\infty} x f_X(x) dx$$

expectation  
measure of additive mean

$$\sigma^2 = D^2[X] := \int_{-\infty}^{+\infty} (x - \xi)^2 f_X(x) dx$$

variance  
measures of dispersion

$$\sigma = D[X]$$

standard deviation

Let  $h: \mathbb{R} \rightarrow \mathbb{R}$  a monotonous function  
of a continuous random variable with density  $f_X$ ,  
then

$$E[h(X)] = \int_{-\infty}^{+\infty} h(x) f_X(x) dx$$

e.g.  $h(x) = e^x$

used later on

if  $X$  normally distributed

then  $Y := e^X$  lognormally distributed

additive mean in this case irrelevant

## 7.1.2 continued

### Multidimensional case

$$X = (X_1, X_2, \dots, X_n)^T : \Omega \rightarrow \mathbb{R}^n$$

n-dimensional random vector

$$\xi = E[X] := (E[X_1], E[X_2], \dots, E[X_n])^T \text{ expectation}$$

Dispersion of two random variables  $X_i$  and  $X_j$

which may be independent from each other, or not

$$\Sigma_X = \text{Cov}[X_i, X_j] := \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2n} \\ \vdots & & \ddots & \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{pmatrix}$$

covariance matrix

variances

$$\sigma_{ij} = \sigma_{ji} = E[(X_i - \xi_i)(X_j - \xi_j)] \text{ covariances}$$

$$\sigma_{ii} = \sigma_i^2 \text{ diagonal elements}$$

$$\exists f \quad \sigma_{ij} = 0 \quad \forall i \neq j$$

then  $X_i$  and  $X_j$  are uncorrelated

Stochastically independent random variables  
are always uncorrelated

(the inverse, however, does not hold, generally)

### 7.1.3 Special Distributions

The Gaussian or normal distribution  $N(\xi, \sigma^2)$

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\xi)^2}{2\sigma^2}}$$

density function  
(bell shaped)

$$F_X(x) = \Phi(x) = \int_{-\infty}^x f_X(x) dx$$

distribution function

with expectation  $\xi$  and standard deviation  $\sigma$

standard normal distribution  $N(0, 1)$

$$\exists f \quad X \sim N(\xi, \sigma^2)$$

$$Y \sim N(0, 1)$$

$$\text{then } Y := \frac{X - \xi}{\sigma} \sim N(0, 1)$$

$$X := \xi + \sigma Y \sim N(\xi, \sigma^2)$$

addition theorem

for stochastically independent  $X_1, X_2, \dots, X_n$   
 $X_i \sim N(\xi_i, \sigma_i^2) \quad \forall i$

$$\sum_{i=1}^n X_i \sim N\left(\sum_{i=1}^n \xi_i, \sum_{i=1}^n \sigma_i^2\right)$$

the sum is also normally distributed!

especially: if  $\xi_i = \xi$  and  $\sigma_i = \sigma \quad \forall i$

$$\text{then } \xi_{\text{sum}} = n\xi, \quad \sigma_{\text{sum}} = \sqrt{n}\sigma$$
$$\sigma_{\text{sum}}^2 = n\sigma^2$$

### 7.1.3 continued

The  $n$ -dimensional normal distribution

$$X = (X_1, X_2, \dots, X_n)^T$$

$$f_X(x) = \frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma_X|}} e^{-\frac{1}{2}(x-\xi)^T \Sigma_X^{-1} (x-\xi)}$$

$\xi$  : vector of expectations

$\Sigma_X$  : covariance matrix ;  $|\Sigma_X|$  its determinant

in this case: stochastic independence

equivalent to uncorrelatedness

off-diagonal elements  $\sigma_{ij} = 0 \quad \forall i \neq j$

$\Sigma_X$  diagonal matrix

The multidimensional normal distribution belongs to the class of elliptical distributions

$Z = RX$  is elliptically distributed if

$X$  is uniformly distributed over an  $n$ -dim. ellipsoid  
and

$R$  is a non-negative random variable that is stochastically independent of  $X$

special case: spherical distribution

if  $X$  is uniformly distributed over a hypersphere

i.e.  $\sigma_i = \sigma \quad \forall i$  and  $\sigma_{ij} = 0 \quad \forall i \neq j$

within an  $n$ -dimensional normal distribution

Then,  $R^2$  is  $\chi^2$  distributed

### 7.1.3 continued II

#### The $\chi^2$ distribution

Let  $X_i \sim N(0, 1) \quad \forall i = 1, 2, \dots, n$

then  $Y := \sum_{i=1}^n X_i^2$  has a  $\chi_n^2$  distribution with  $n$  degrees of freedom

$$E[Y] = n$$

$$D^2[Y] = 2n$$

$\exists$  if  $n \gg 1$  the  $\chi^2$  distribution can be approximated by an  $N(n, 2n)$  distribution

i.e.  $\chi_n^2 \rightarrow n + \sqrt{2n} N$  for  $n \gg 1$

$R := +\sqrt{Y}$  may be handled as if  $R \sim N(\sqrt{n}, \sqrt{\frac{1}{2}})$

(follows from rules for  $E[h(x)]$   
 $D^2[h(x)]$  above)

### 7.1.3 continued III

#### Order Statistics

Let  $X_1, X_2, \dots, X_n$  be stochastically independent random variables with  $F_{X_i}$  and  $f_{X_i}$

then  $Y_1 = \max \{X_1, X_2, \dots, X_n\}$

the distribution function  $F_{Y_1}$  of which is

$$\begin{aligned} F_{Y_1}(x) &= \prod_{i=1}^n F_{X_i}(x) \\ &= p(X_1 \leq x) \cdot p(X_2 \leq x) \cdot \dots \cdot p(X_n \leq x) \end{aligned}$$

$\exists$  if all  $X_i$  are identically distributed

$$F_{Y_1}(x) = [F_X(x)]^n$$

$$f_{Y_1} = \frac{\partial}{\partial x} F_{Y_1} = n f_X(x) [F_X(x)]^{n-1}$$

More generally:

Let  $Y_m$  be the  $m$ -th largest  $X_i$

then

$$f_{Y_m}(x) = \underbrace{n \binom{n-1}{m-1}}_c f_X(x) \underbrace{[F_X(x)]^{n-m}}_b \underbrace{[1 - F_X(x)]^{m-1}}_a$$

a:  $m-1$  larger

b:  $n-m$  smaller

elements  $X_i$  in vector  $X$

c: combinatorial manifold of situations



## 7.2 Convergence behavior of ESs

|      |  |   |  |                               |
|------|--|---|--|-------------------------------|
| 1951 | $(1+1)$<br>$(\mu+1)$   | two membered<br>first multimembered                   |  | sequential<br>Rechenberg      |
| 1974 | $(1+\lambda)$ , $(1, \lambda)$<br>$(\mu+\lambda)$ , $(\mu, \lambda)$ |   |  | parallel<br>Schwefel          |
|      |  | comma-versions: non-elitist<br>plus-versions: elitist |  | $\kappa=1$<br>$\kappa=\infty$ |
| 1995 | $(\mu, \kappa, \lambda)$ ES  | $\kappa = \text{max. life span}$                      |  |                               |

generally: theory can assure something only on the basis on certain assumptions (simplifications)

efficiency results need more rigorous ass. than effectiveness results

nevertheless:

except by chance

no method can be better  $\checkmark$  in the more general case than has been proven for simplified situations

e.g. Newton-Raphson method

just one step necessary if objective quadratic

but may be divergent in slightly modified situations

EAs aim at robustness

however: no theory for general efficiency

## 7.2.1 Convergence reliability

here: only continuous case  $x \in \mathbb{R}^n$

[see Rudolph, 1994 for  $x \in \mathbb{N}^n$ ]

here: only  $(1+1)ES$

[see Rudolph, 1995 for comma-versions]

We start at  $x^{(0)}$  and iterate

$$x^{(k+1)} = x^{(k)} + z \quad \text{with } z \sim N(0, \sigma^2 I_n)$$

We assume a regular optimization problem, i.e.

$$f^* = f(x^*) = \min \{ f(x) \mid x \in M \subseteq \mathbb{R}^n \}$$

a) with  $f^* > -\infty$

b)  $x^* \in \text{int}(M)$

interior

c)  $\mu(\{x \in M \mid f(x) \in U_\varepsilon(f^*)\}) > 0 \quad \forall \varepsilon > 0$

Lebesgue measure

[ $f^*$  global minimum;  $x^*$  solution or global minimizer]

- a) necessary, otherwise there is no finite minimizer
- b) only simplifies analysis
- c) excludes singular optima (needle in a haystack) that cannot be found with probability  $> 0$

## 7.2.1 continued

### Theorem A

Let  $\varepsilon > 0$  and  $p_t := p(x^{(t)} \in \{x \in M \mid f(x) \in U_\varepsilon(f^*)\})$   
the prob. that a  $(1+\varepsilon)$ ES has reached  
 $x^{(t)}$  at iteration  $t$   
the obj. function value of which  
is closer to  $f^*$  than  $\varepsilon$  ;

then assuming that

$$\sum_{t=0}^{\infty} p_t = \infty$$

A<sub>1</sub>

implies

$$p(\lim_{t \rightarrow \infty} (f(x^{(t)}) - f^*) = 0) = 1$$

for any starting point  $x^{(0)} \in M$

more practically (Lemma A') :

$\exists f$   $M \subseteq \text{support}(f_z)$  with  $f_z$  probab. dens. of vector  
 $z$  of the mutation operator, and  $M$  is bounded, then  
A<sub>1</sub> is valid

more loosely:

you can reach  $x^*$  at any iteration  $t$  immediately  
as long as  $\sigma^{(t)} > 0$

This result is of academic interest only!

More interesting: Time to achieve a certain  
approximation

## 7.2.2 Convergence velocity (or rate)

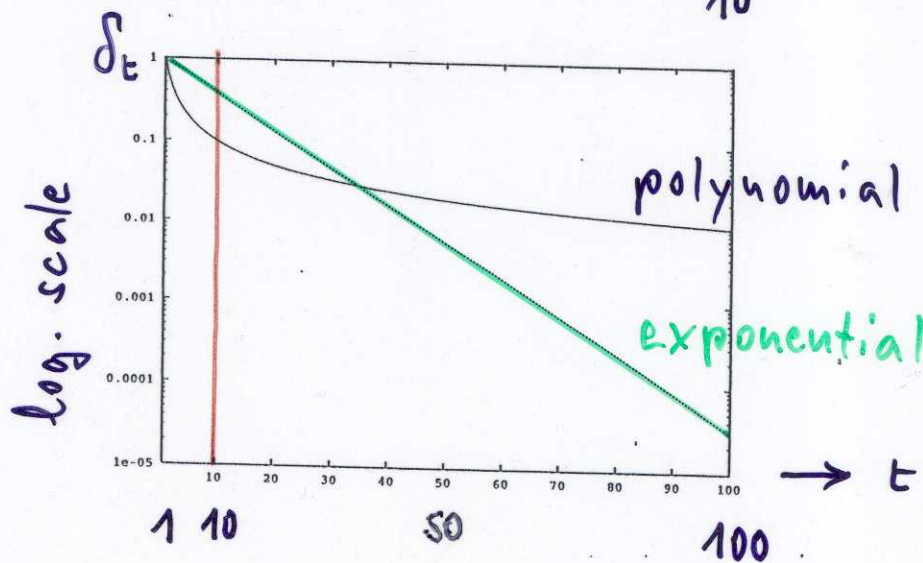
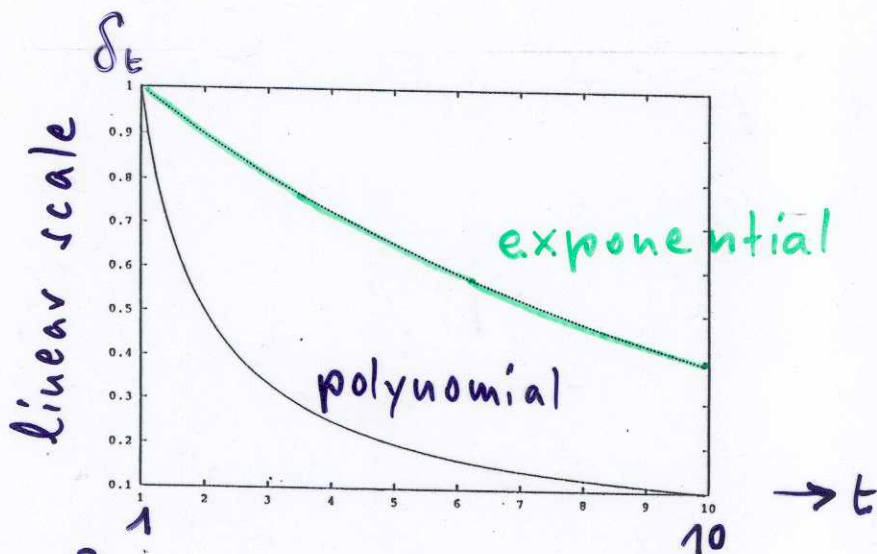
Let  $\delta_t := E[f(x^{(t)}) - f^*]$  expected error at step  $t$

There are different kinds of convergence order

$$\delta_t = \mathcal{O}(t^{-\alpha}), \quad \alpha > 0 \quad \text{polynomial convergence order}$$

$$\delta_t = \mathcal{O}(\beta^t), \quad \beta \in (0, 1) \quad \text{exponential}$$

→ two graphs



## 7.2.2 continued

special assumptions necessary

here: strictly convex problems

(implies continuous differentiability, unimodality)

This is not the domain of EAs, but

any method that is good for difficult problems should behave even better in simple situations.

Thus, we may gain best-case results, now.

(stronger ones in the future, hopefully)

### Theorem B

Let  $f: M \rightarrow \mathbb{R}$  strictly convex objective function

mut. step size of a  $(1+1)$  ES spherically distrib.

with  $Z = \mathbb{R}^n$ , support  $(R) = (0, a) \neq \emptyset$

Then 
$$\delta_t = \begin{cases} \Theta(t^{-2/n}) & \text{for a constant mut. step size} \\ \Theta(\beta^t) & \text{for an adaptive} \end{cases}$$

at step  $t$  for any start position  $x^{(0)} \in M$

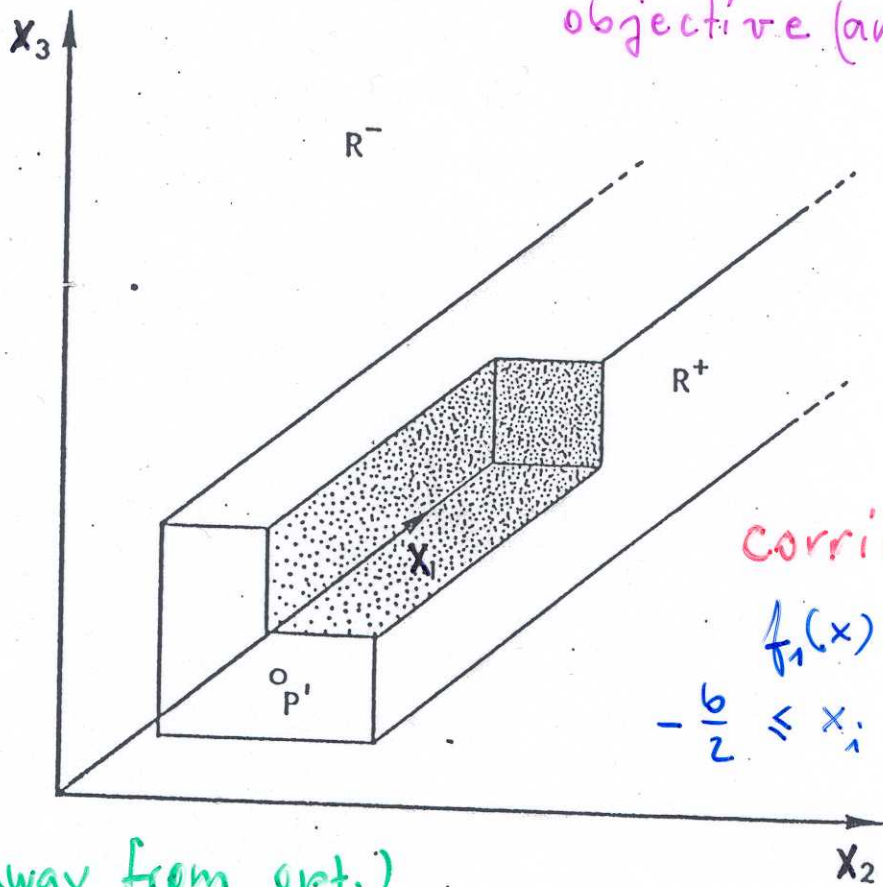
with  $\beta \in (0, 1)$  and step size adaptation

according to  $R^{(t+1)} = \|\nabla f(x^{(t)})\| R^{(t)}$ .

or  $2/5$  success rule (Rechenberg, Rapp1)

How large is  $\beta$ ?

two (n-dimensional) fitness functions  
 objective (and constraints)

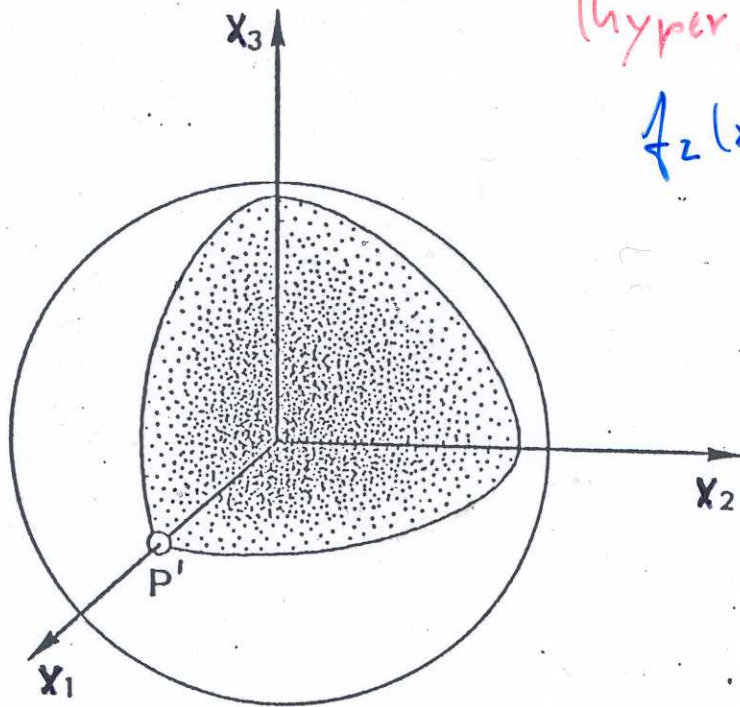


corridor (ridge)

$$f_1(x) = -a x_1$$

$$-\frac{b}{2} \leq x_i \leq \frac{b}{2} \quad \forall i=2,3,\dots,n$$

(far away from opt.)



(hyper)-sphere

$$f_2(x) = \sum_{i=1}^n (x_i - x_i^*)^2$$

$$= r^2$$

(close to opt.)

i.e. two different 'typical' situations

## 7.3 The Sphere Model n-dimensional

$$f(x) := \|x - x^*\|^2 = r^2$$

Results are valid for many more models with concentric hypersurfaces  $f(x) = \text{const.}$

What we want to know:  $n$ -dependency and  $\sigma$ -

of  $\varphi := E[r^{(t)} - r^{(t+1)}]$   
mean convergence velocity  $\varphi(n, \sigma, r)$

for ESs with just one common  $\sigma$

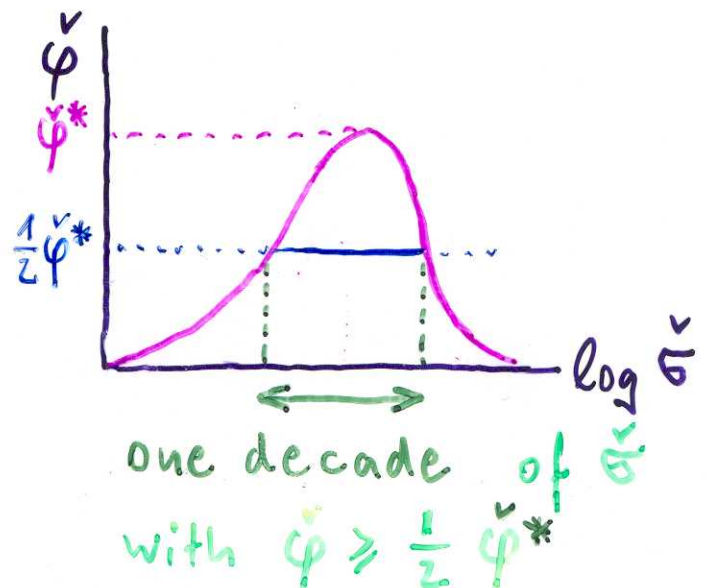
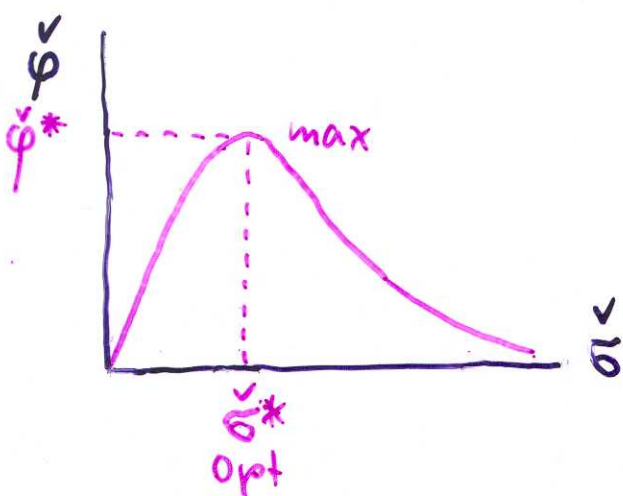
most interesting:  $\varphi_{\max}(n, \sigma_{\text{opt}}, r)$

### 7.3.1 The (1+1) ES / two membered

$$\varphi_{1+1} \approx \frac{r}{n} \left\{ \frac{\sigma n}{\sqrt{2\pi} r} e^{-\left(\frac{\sigma n}{\sqrt{2} r}\right)^2} - \left(\frac{\sigma n}{\sqrt{2} r}\right)^2 \left[1 - \Phi\left(\frac{\sigma n}{\sqrt{2} r}\right)\right] \right\}$$

result

with 'dimensionless'  $\check{\varphi} = \varphi \frac{n}{r}$ ,  $\check{\sigma} = \sigma \frac{n}{r}$  for  $n \gg 1$



in between:

$$\begin{aligned}
 r^{(t+1)} &= \|x^{(t)} - x^* + z\| = \sqrt{\sum_{i=1}^n ((x_i^{(t)} - x_i^*) + \sigma N_i)^2} \\
 &= \sqrt{\sum_{i=1}^n (x_i^{(t)} - x_i^*)^2 + 2\sigma \sum_{i=1}^n (x_i^{(t)} - x_i^*) N_i + \sigma^2 \sum_{i=1}^n N_i^2} \\
 &\sim \sqrt{r^2 + 2\sigma r N + \sigma^2 (n + \sqrt{2n} N')} \leftarrow \underbrace{\sigma^2 \sum_{i=1}^n N_i^2}_{\chi^2} \\
 &\sim \sqrt{r^2 + n\sigma^2 + \sqrt{2n\sigma^4 + 4r^2\sigma^2} N''}
 \end{aligned}$$

addition th.

$P := \frac{r^{(t+1)} - r}{r}$  relative progress

$$P = 1 - \sqrt{1 + \frac{n\sigma^2}{r^2} + \sqrt{\frac{2n\sigma^4}{r^4} + 4\frac{\sigma^2}{r^2} N''}}$$

with  $\sigma \frac{n}{r} = \sigma^{\check{}}$

$$P = 1 - \sqrt{1 + \frac{\sigma^{\check{2}}}{5} + \frac{\sigma^{\check{2}}}{5} \sqrt{\frac{4}{\sigma^{\check{2}}} + \frac{2}{n} N''}}$$

small if n large

$$\approx 1 - \sqrt{1 + \frac{\sigma^{\check{2}}}{5} + \frac{2\sigma^{\check{2}}}{5} N''}$$

with  $\sqrt{1+\epsilon} \approx 1 + \frac{\epsilon}{2}$

$$\approx 1 - \left(1 + \frac{\sigma^{\check{2}}}{2n} + \frac{\sigma^{\check{2}}}{5} N''\right) = -\frac{\sigma^{\check{2}}}{2n} - \frac{\sigma^{\check{2}}}{5} N''$$

$$\sim N\left(-\frac{\sigma^{\check{2}}}{2n}, \frac{\sigma^{\check{2}}}{5}\right)$$

$$E[P_{t+1}] = \int_0^{\infty} x f_P(x) dx := \frac{\varphi_{t+1}}{r} = \frac{\check{\varphi}_{t+1}}{n}$$

only positive moves accepted  
within  $(t+1)\epsilon$

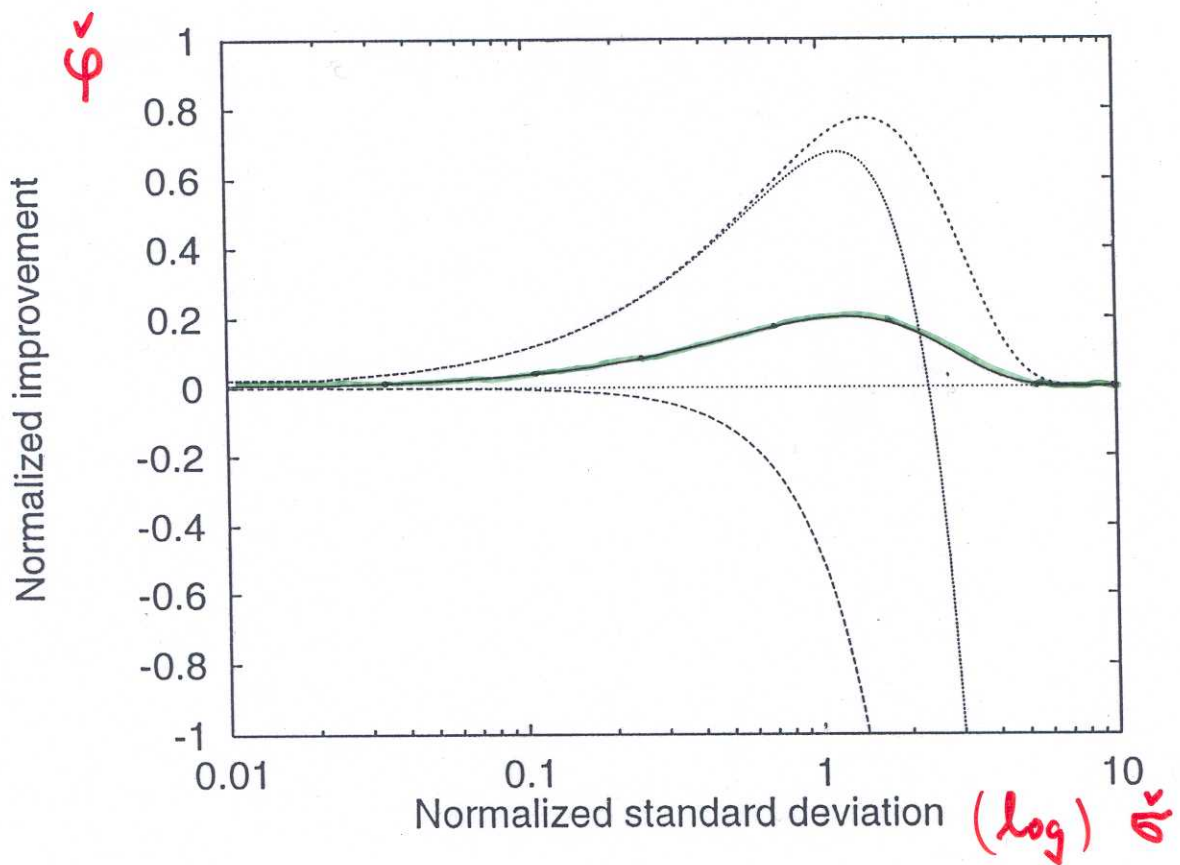
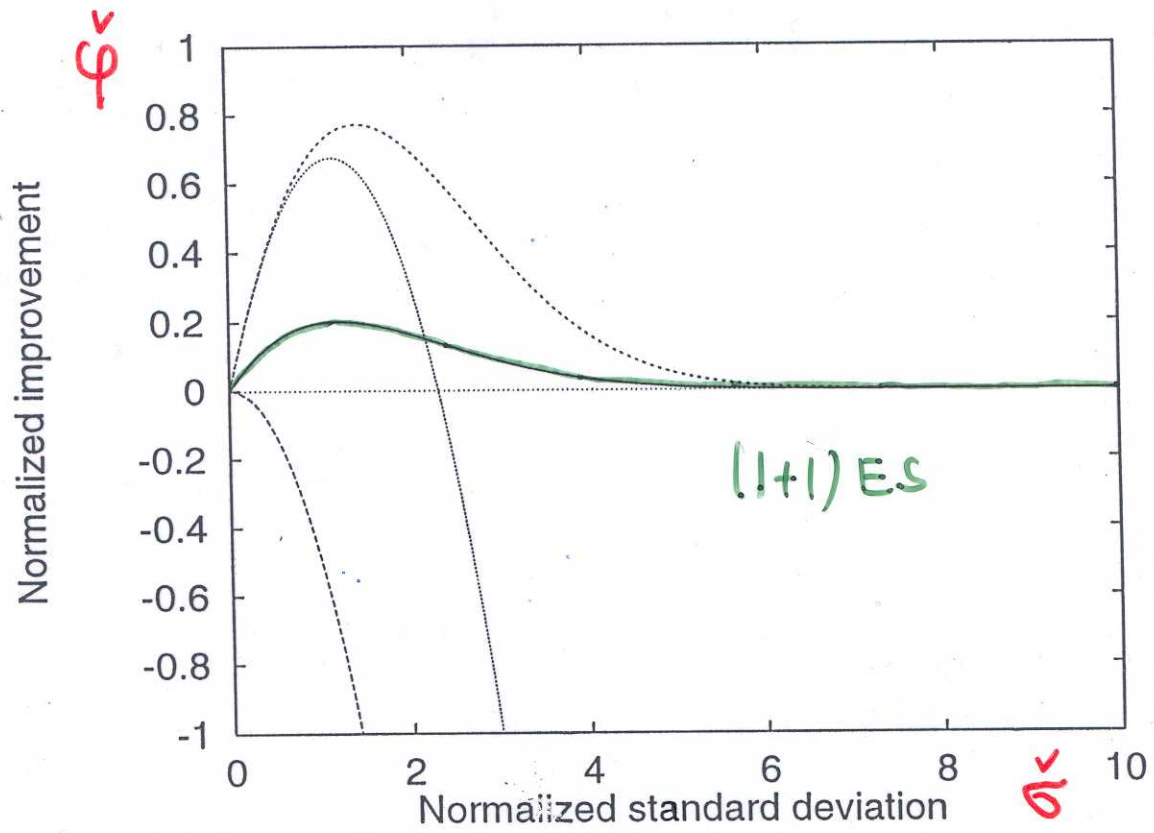


$$f_p(x) = \frac{1}{\sqrt{2\pi} \frac{\sigma}{5}} e^{-\frac{1}{2} \left( \frac{x + \frac{\sigma}{5}}{\frac{\sigma}{5}} \right)^2}$$

$$= \frac{n}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \left( \frac{nx}{\sigma} + \frac{\sigma}{2} \right)^2}$$

$$E[P_{i+1}] = \int_0^{\infty} \frac{nx}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \left( \frac{nx}{\sigma} + \frac{\sigma}{2} \right)^2} dx$$

$$= \frac{1}{5} \left[ \frac{\sigma}{\sqrt{2\pi}} e^{-\frac{\sigma^2}{2}} - \frac{\sigma^2}{2} \left[ 1 - \Phi\left(\frac{\sigma}{2}\right) \right] \right]$$



### 7.3.1 continued

$$\varphi_{l+1}^* = 0.2025 \frac{\tau}{n} \quad \text{at} \quad \sigma^* = 1.224 \frac{\tau}{n} = \frac{0.612}{n} \|\nabla f(x)\|$$

$$\sim \frac{1}{n}$$

$$\sim \frac{1}{n}$$

inversely proportional  
to the # of variables

Success probability

$$p_s^* \approx \Phi\left(-\frac{\sigma^*}{2}\right) \approx 0.270$$

Similar results available for an  $n$ -dimensional  
'corridor' model of width  $b$

$$\varphi_{l+1}^* \sim \frac{b}{n}, \quad \sigma^* \sim \frac{b}{n}, \quad p_s^* = \frac{1}{2e} \approx 0.184$$

⇒ formulation of 1/s success rule

### 7.3.2 The $(\mu+1)$ ES

extinction of the worst

Similar to EVOP strategy of G.E.P. Box

polyhedron strategies

Simplex

Nelder and Mead

Complex

M. J. Box

no longer in use

since self-adaptation of  $\sigma$  does not work  
within the  $(\mu+1)$  scheme

### 7.3.3 The $(1+\lambda)$ ES

Already calculated: relative change per generation of each successor

$$\rho \approx -\frac{\check{\sigma}^2}{2n} + \frac{\check{\sigma}}{n} N''$$

The corresponding distribution function is

$$F_{\rho}(x) = \Phi\left(\frac{x-\theta}{\eta}\right)$$

with expectation  $\theta = -\frac{\check{\sigma}^2}{2n}$

and standard deviation  $\eta = \frac{\check{\sigma}}{n}$

From order statistics:

The best out of  $\lambda$  successors

$$F_{\rho_{\lambda}}(x) = \Phi^{\lambda}\left(\frac{x-\theta}{\eta}\right)$$

with  $f_{\rho_{\lambda}}(\hat{x}) = \frac{d}{dx} \Phi^{\lambda}(\hat{x}) = \frac{\lambda}{\sqrt{2\pi}} e^{-\frac{\hat{x}^2}{2}} \Phi^{\lambda-1}(\hat{x})$

from above

$$E[\rho_{\lambda}] = \frac{1}{\eta} \int_{-\infty}^{\infty} x f_{\rho_{\lambda}}\left(\frac{x-\theta}{\eta}\right) dx$$

$$E[\rho_{1+\lambda}] = E[\max\{0, \rho_{\lambda}\}]$$

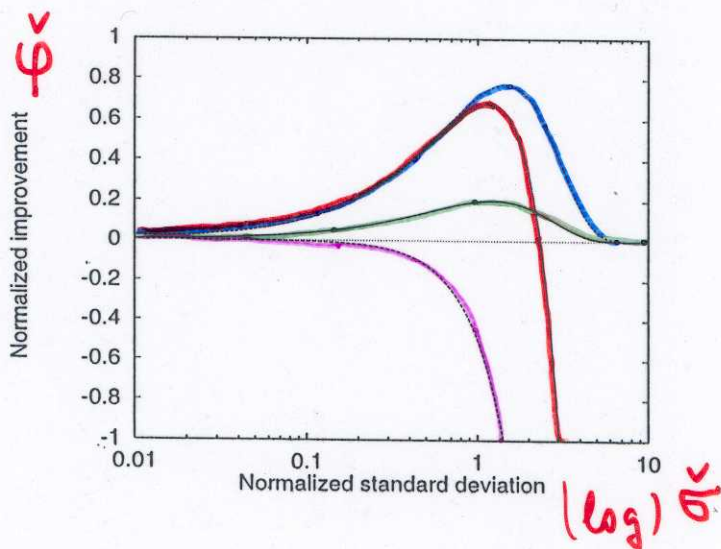
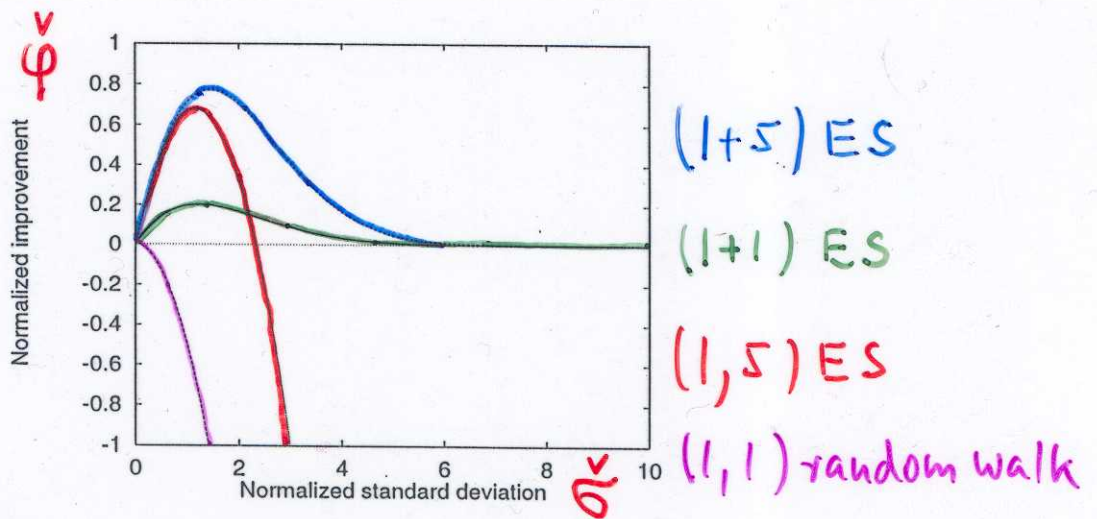
$$= \frac{1}{\eta} \int_0^{\infty} x f_{\rho_{\lambda}}\left(\frac{x-\theta}{\eta}\right) dx$$

Substituting  $z := (x - \theta) / \eta$  yields

$$E[V_{1+\lambda}] = \int_{-\theta/\eta}^{\infty} (z\eta + \theta) f_{P_\lambda}(z) dz$$

$$= \frac{\sigma^2}{\eta} \int_{\sigma^2/2}^{\infty} z f_{P_\lambda}(z) dz - \frac{\sigma^2}{2\eta} [1 - \Phi^\lambda(\frac{\sigma^2}{2})]$$

integral must be evaluated numerically



### 7.3.4 The $(1, \lambda)$ ES

non-elitist!  $\Rightarrow$  it may diverge if  $\sigma$  too large

Rudolph (1994) has shown: under certain conditions the  $(1, \lambda)$  ES converges on average almost sure.

For mean convergence rate we just have to change the lower bound of the integral only

$$\begin{aligned} E[P_{1,\lambda}] &= \int_{-\infty}^{\infty} (z\eta + \theta) f_{P_{\lambda}}(z) dz \\ &= \frac{1}{n} \underbrace{\int_{-\infty}^{\infty} z f_{P_{\lambda}}(z) dz}_{c_{1,\lambda} \text{ new constant}} - \frac{\check{\sigma}^2}{2n} \\ &= (2\check{\sigma} c_{1,\lambda} - \check{\sigma}^2) / 2n \end{aligned}$$

Looking for a maximum at some optimal  $\check{\sigma}^*$

$$\frac{\partial}{\partial \check{\sigma}} E[P_{1,\lambda}] = \frac{1}{2n} (2c_{1,\lambda} - 2\check{\sigma})$$

becomes zero if  $\check{\sigma} = c_{1,\lambda}$

$$\frac{\partial^2}{\partial \check{\sigma}^2} E[P_{1,\lambda}] = -\frac{1}{n} < 0$$

Thus, the optimal mutation step size is

$$\sigma^* = c_{1,\lambda} \frac{\tau}{n} = \frac{c_{1,\lambda}}{2n} \|\nabla f(x)\|$$

$$E^*[P_{1,\lambda}] = \frac{1}{2n} c_{1,\lambda}^2$$

### 7.3.4 continued

in other terms

$$\check{\varphi}_{1,\lambda}^* = \frac{1}{2} c_{1,\lambda}^2 = \frac{1}{2} \check{\sigma}^{*2}$$

$c_{1,\lambda}$ : tabulated

For large  $n$  and large  $\lambda$  it has been shown that

$$c_{1,\lambda} \sim \sqrt{2 \ln \lambda}$$

$$\check{\varphi}_{1,\lambda}^* \sim \ln \lambda$$

The same holds for the  $(1+\lambda)ES$ !

$$\check{\varphi}_{1+\lambda}^* \sim \ln \lambda$$

⇒ logarithmic speedup

if  $\lambda$  successors are evaluated in parallel

### 7.3.5 The $(\mu, \lambda)$ ES

Beyer (1994) has shown in a rather lengthy paper:

$$\varphi_{\mu, \lambda}^* \sim \ln \frac{\lambda}{\mu}$$

$\Rightarrow$  It's the selection pressure that speeds up

Is  $\mu = 1$  optimal? No!

One may even achieve linear speedup by increasing  $\mu$  at constant  $\lambda/\mu$

with recombination!

Beyer (1994) investigated panmictic recomb.

- intermediary (globally)
- global discrete

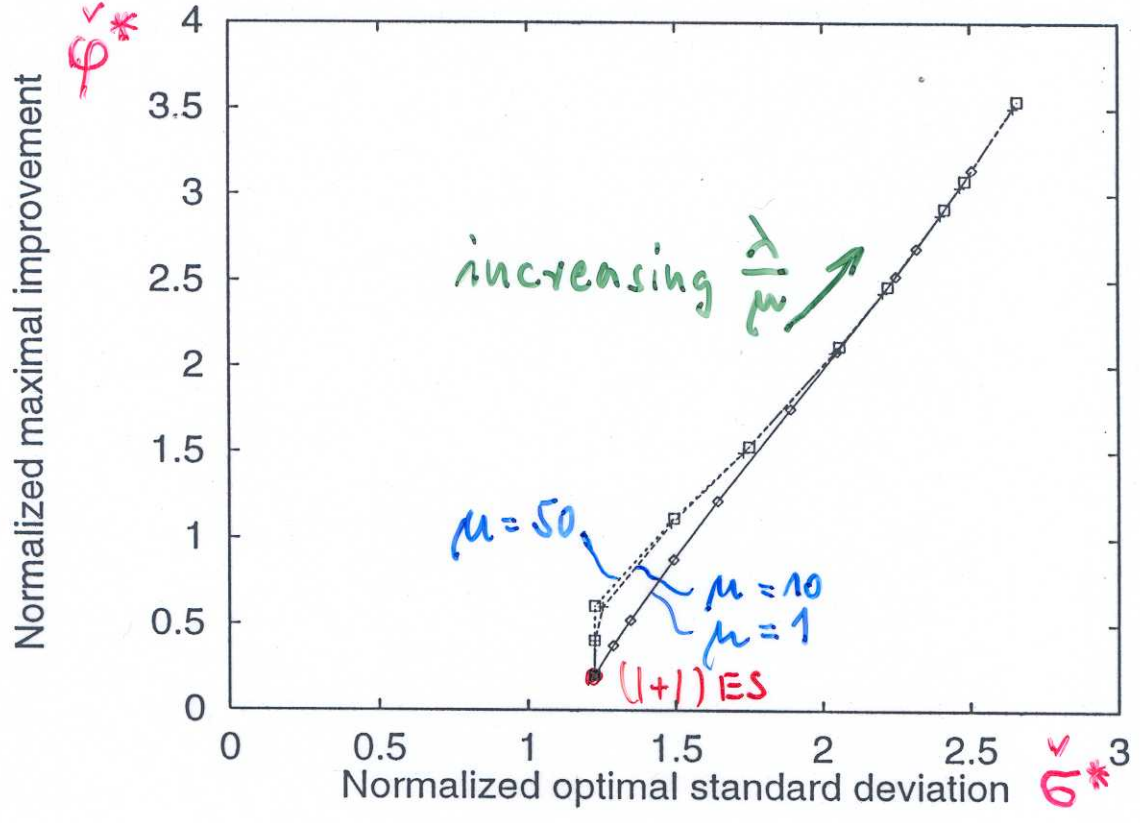
I add a bit speculation (supported by graphs):

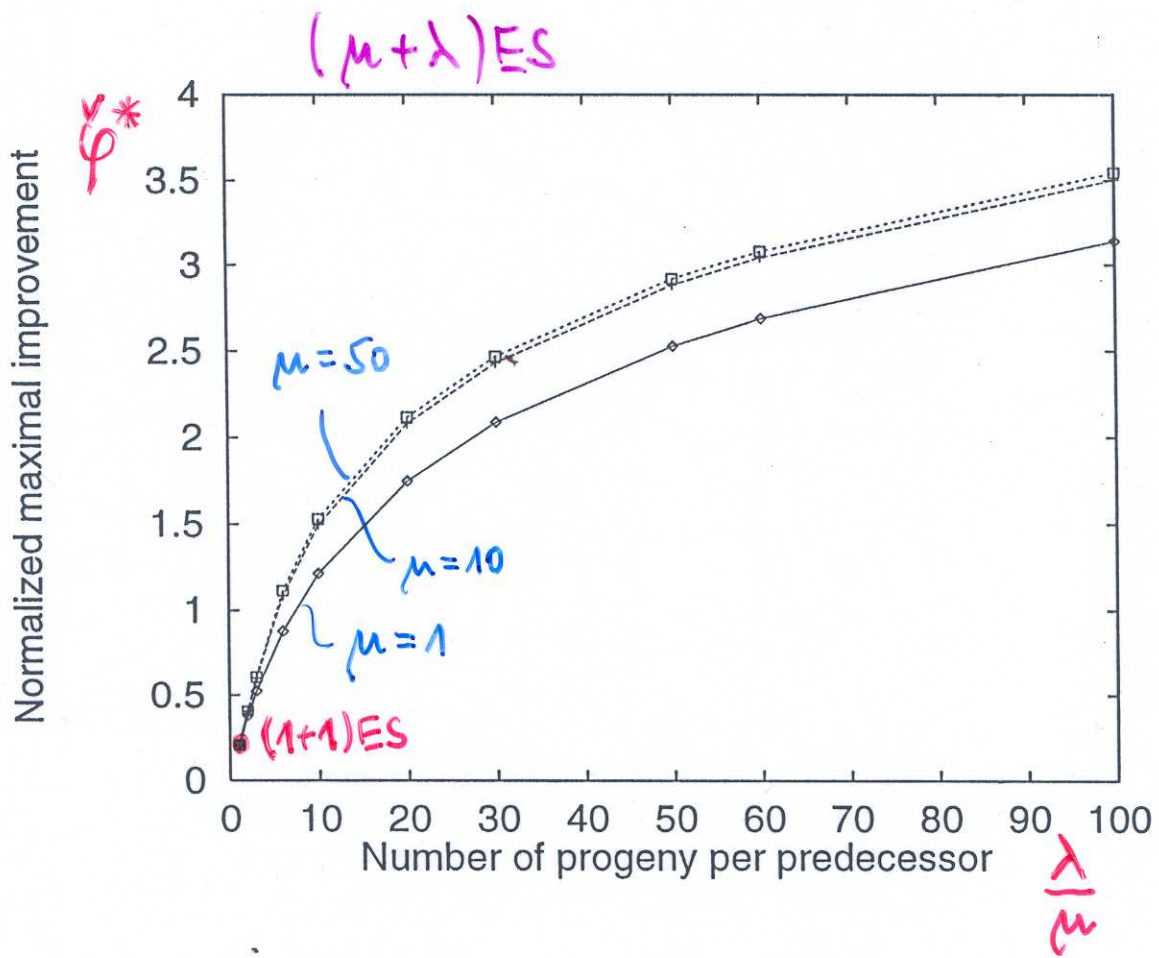
$$\varphi_{\mu, \kappa=1, \lambda, \sigma=\mu}^* \sim \mu \ln \frac{\lambda}{\mu}$$

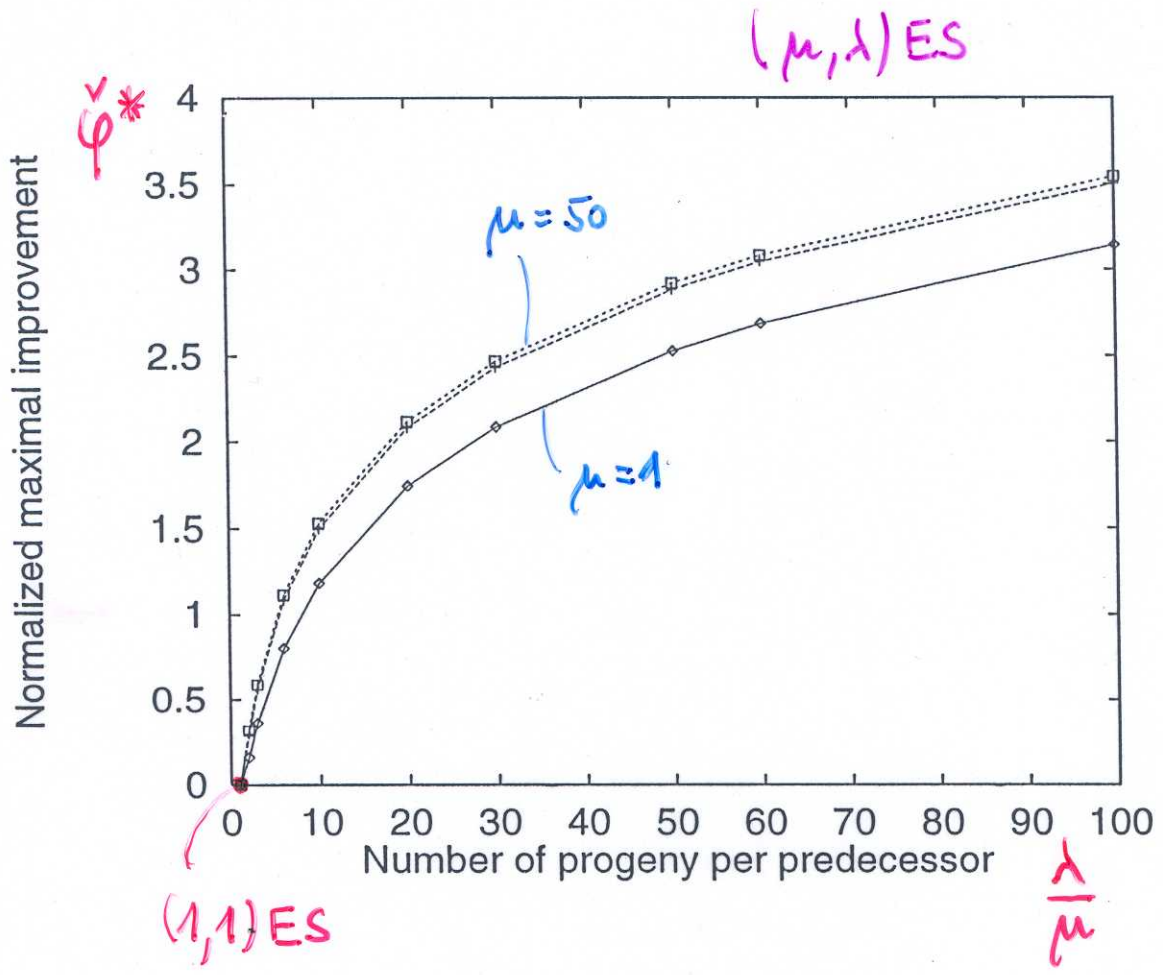
$$\sigma^* = \begin{cases} \sqrt{2\varphi^*} & \text{for uniform crossover} \\ \sqrt{2\mu\varphi^*} & \text{for global interm. reco.} \end{cases}$$

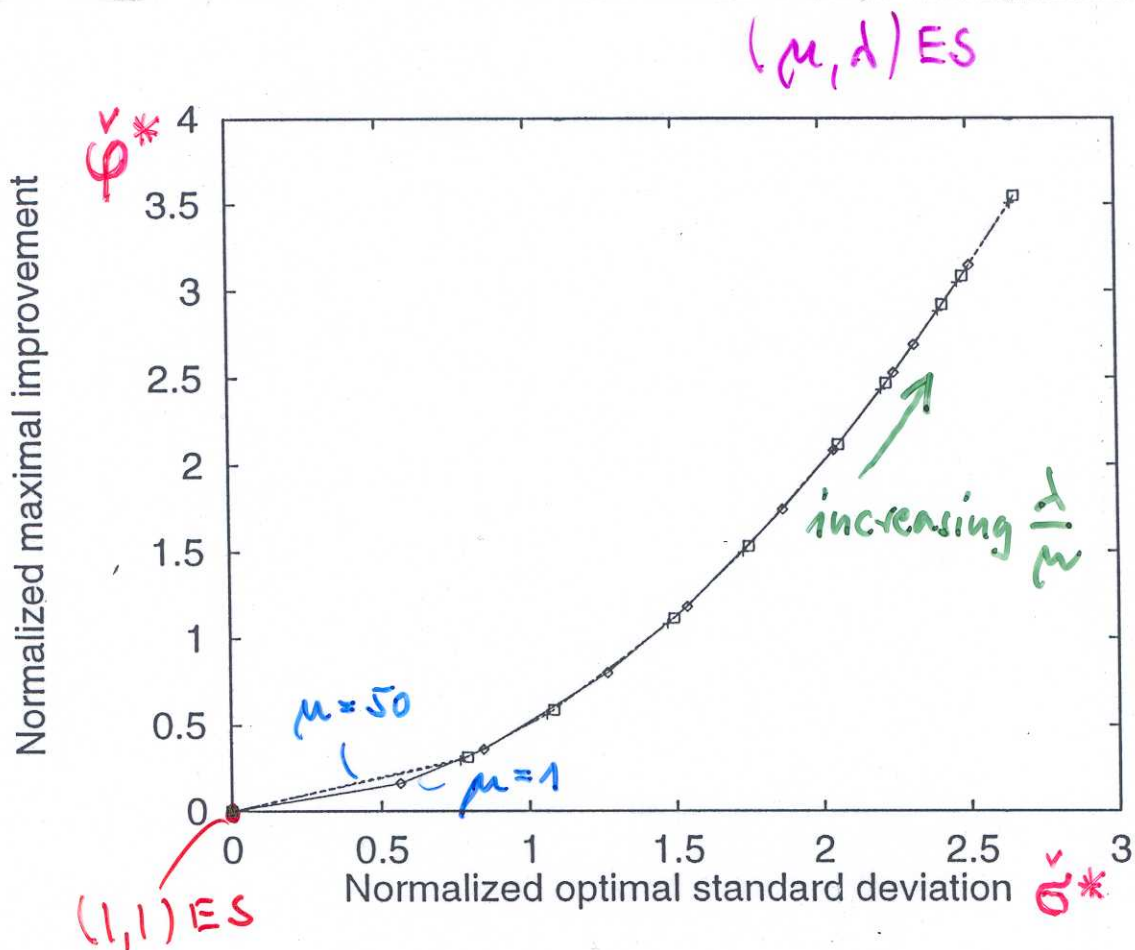


$(\mu + \lambda) ES$









conclusion :

efficiency mainly depends on selection pressure  
(if  $\mu$  and  $\lambda$  not too small)

without recombination :

population size not thus important  
(if  $\mu$  and  $\lambda$  not too small)

if  $\mu$  and  $\lambda$  not too small and  $\frac{\lambda}{\mu} > 6$  :  
no difference between

comma-ES and  
plus-ES (elitist version)