



HAL
open science

Stagnation Detection meets Fast Mutation

Benjamin Doerr, Amirhossein Rajabi

► **To cite this version:**

Benjamin Doerr, Amirhossein Rajabi. Stagnation Detection meets Fast Mutation. Evolutionary Computation in Combinatorial Optimization (EvoCOP 2022), Apr 2022, Madrid, Spain. 10.1007/978-3-031-04148-8_13 . hal-03797608v1

HAL Id: hal-03797608

<https://hal.science/hal-03797608v1>

Submitted on 4 Oct 2022 (v1), last revised 30 Mar 2023 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Stagnation Detection meets Fast Mutation

Benjamin Doerr* Amirhossein Rajabi†

January 31, 2022

Abstract

Two mechanisms have recently been proposed that can significantly speed up finding distant improving solutions via mutation, namely using a random mutation rate drawn from a heavy-tailed distribution (“fast mutation”, Doerr et al. (2017)) and increasing the mutation strength based on stagnation detection (Rajabi and Witt (2020)). Whereas the latter can obtain the asymptotically best probability of finding a single desired solution in a given distance, the former is more robust and performs much better when many improving solutions in some distance exist.

In this work, we propose a mutation strategy that combines ideas of both mechanisms. We show that it can also obtain the best possible probability of finding a single distant solution. However, when several improving solutions exist, it can outperform both the stagnation-detection approach and fast mutation. The new operator is more than an interleaving of the two previous mechanisms and it also outperforms any such interleaving.

1 Introduction

Leaving local optima is a challenge for evolutionary algorithms. Mutation-based approaches are challenged by the fact that the typical mutation rate of $p = 1/n$ rarely leads to offspring in a larger distance from the parent. When using larger mutation rates, the choice of the mutation rate is critical and small constant-factor deviations from the optimal rate can lead to huge performance losses [DLMN17, Cor. 4.2].

*Laboratoire d’Informatique (LIX), CNRS, École Polytechnique, Institut Polytechnique de Paris, Palaiseau, France

†Technical University of Denmark, Kgs. Lyngby, Denmark

Two ways to overcome this problem were proposed recently, namely the use of a random mutation rate sampled from a power-law distribution [DLMN17] and the successive increase of the mutation rate when a stagnation-detection mechanism indicates that the current rate is unlikely to generate solutions not seen yet [RW20]. An improved version of this stagnation-detection approach [RW21b], the so-called SD-RLS algorithm based on k -bit mutation instead of standard bit mutation, can find a single improving solution in distance m in expected time $(1 + o(1))\binom{n}{m}$ (without knowing that the distance to the desired solution is m). Apart from lower order terms, this is the same runtime that can be obtained via a repeated use of the best unbiased mutation operator that is aware of m (which is, naturally, flipping m random bits). It is faster than the fast $(1 + 1)$ EA by a factor of $\Omega(m)$.

While the SD-RLS algorithm thus is very efficient in finding a single desired solution (and thus has very good runtimes on the classic jump functions benchmark), this algorithm has a poor performance when there are several improving solutions in distance m as now the stagnation detection approach leads to too much time spent on too small mutation strength. Taking the generalized jump function [BBD21] having a valley of low fitness of width δ , $\delta \geq 2$ a constant, in distance $n/4$ from the optimum as extreme example, we easily see that the SD-RLS takes an expected time of $\Omega(n^{\delta-1})$ to traverse the fitness valley, whereas the $(1 + 1)$ EA both with the classic mutation operator and with fast mutation does so in expected constant time.

Our results: Based on this insight that fast mutation and stagnation detection have complementary strengths, we design a mutation-based approach that takes inspiration from both approaches. We follow, in principle, the basic version of the improved stagnation-detection approach of [RW21b], that is, we start with mutation strength $r = 1$ and increase r gradually. More precisely, when strength r has been used for a certain number ℓ_r of iterations without that an improvement was found, we increase r by one since we assume that no improvement in distance r exists (we omit some technical details in this first presentation of our approach, e.g., that we do not increase r beyond $n/2.1$, and refer the reader to Algorithm 1 for the full details). Different from [RW21b], when the current strength is r , we do not always flip r random bits as mutation operator, but we choose a random number X_r of bits to flip. This number is equal to r with probability $1 - \gamma$, where γ is an algorithm parameter that is usually small (a small constant or $o(1)$). With probability γ , however, X_r deviates from r by an amount following a power-law distribution with exponent β . The precise definition of this case (see again Algorithm 1) is not too important, so for this first exposition we can assume that we sample D from a power-law distribution (with exponent

β) on the positive integers and then, each with probability $1/2$, flip $r + D$ or $r - D$ random bits (where we do nothing if this number is not between 1 and n).

Since with probability $1 - \gamma$, we essentially follow the basic approach of [RW21b], it is not surprising that we find a single closest improving solution in distance m in an expected time of $\frac{1}{1-\gamma}(1 + o(1))\binom{n}{m}$, again without that the algorithm needs to know m (Theorem 5). If $\gamma = o(1)$, this is again the optimal time of $(1 + o(1))\binom{n}{m}$ discussed above. We note, however, that our algorithm is simpler than the solution presented in [RW21b]. The basic SD-RLS algorithm proposed in [RW21b] obtains a runtime of $(1 + o(1))\binom{n}{m}$ only with high probability and otherwise fails. To turn this algorithm into one that never fails and has an expected runtime of $(1 + o(1))\binom{n}{m}$, a robust version of the SD-RLS was developed in [RW21b] as well. This version repeats previous phases as follows. When the ℓ_r uses of strength r have not led to an improvement, before increasing the rate to $r + 1$, first another ℓ_i iterations are performed with strength i , for $i = r - 1, \dots, 1$. In our approach, such an additional effort is not necessary since the fast mutations automatically render the algorithm robust.

The use of a heavy-tailed mutation rate also helps in situations where the stagnation-detection mechanism takes too long to use larger mutation strengths. Since in phases $r = 1, \dots, 2m$ the probability to flip m bits is at least $\gamma/2$ times the probability of this event in a run of the fast $(1 + 1)$ EA, it is not surprising that our algorithm finds an improvement in distance m is at most $2/\gamma$ times the time of the fast $(1 + 1)$ EA, which as discussed above can be significantly faster than the SD-RLS. Such a result could also have been obtained from a simple interleaving of SD-RLS and fast $(1 + 1)$ EA iterations. Since our heavy-tailed choices of the mutation strength, however, take into account the current strength r , we often obtain better runtimes, often better than both the SD-RLS and the fast $(1 + 1)$ EA. Since the precise statement of these results is technical, we defer the details to Section 4. As a simple example showing the outperformance of our algorithm, we regard the generalized jump function $\text{JUMP}_{m,\delta:=m-\Delta}$ for a constant value of $\Delta \geq 2$ and $m = \omega(1)$. This jump function is similar to the classic jump function JUMP_m , but the valley of low fitness consists not of all search points in positive distance at most $m - 1$ from the optimum, but only of those in distance $\Delta + 1, \dots, m - 1$. Consequently, from the local optimum there is not a single improving solution, but $\Theta(n^\Delta)$. Note that this is still relatively few compared to the fitness valley of size essentially $\binom{n}{m}$. On this generalized jump function, the expected runtimes of SD-RLS is $O\left(\binom{n}{\delta-1} \ln(R)\right)$, the one of the fast $(1 + 1)$ EA is $O(\delta^{\beta-0.5}(en/\delta)^\delta n^{-\Delta})$, and the one of our algorithm

is at most $O\left(\binom{n}{\delta} n^{-\Delta} \gamma^{-1}\right)$ (Corollary 11). Since it is also clear that any interleaving of SD-RLS and fast $(1 + 1)$ EA iterations cannot give a better runtime than the one of the two pure algorithms, this result shows that our algorithm can beat SD-RLS and fast EA (and any simple mix of them) when there are several improving solutions in a given distance.

Structure of this paper: After reviewing the most relevant previous works in Section 2, we introduce our new algorithm in Section 3. In Section 4, we analyze via mathematical means how our algorithm finds an improvement in distance m both when this is typically achieved in phase m (e.g., when there is only one improving solution in distance m) and when this is achieved earlier via the heavy-tailed rates. We use these results in Section 5 to prove several runtime results, among others, for generalized jump functions. We present some experimental results in Section 6. In Section 7, we discuss recommendations on how to set the parameters of our algorithm. We conclude the paper with a short discussion of our results and a pointer to possible future work in Section 8.

2 Previous Works

This work aims at combining the advantages of stagnation detection and heavy-tailed mutation, so clearly these topics contain the most relevant previous works. Both integrate into the wider questions of how to optimally set the mutation strength of evolutionary algorithms (for this we refer to the recent survey [DD20]) and how evolutionary algorithms can leave local optima (here we refer to [Doe20, Section 2.1] for a discussion of non-elitist approaches and to the introduction of [DFK⁺18] for a discussion of crossover-based approaches).

For elitist mutation-based approaches, it is clear that when the population has converged to a local optimum the only way to leave this is by mutating a solution from the local optimum into an at least as good solution outside this local optimum. It was observed in [DLMN17] (the earlier work [Prü04] contains similar findings for the special case that the nearest improving solution is in Hamming distance two or three) that standard-bit mutation with mutation rate $p = \frac{1}{n}$, which is the most recommended way of doing mutation, is not perfectly suitable to perform larger jumps in the search space. In fact, when the nearest improving solution is in Hamming distance m , then a mutation rate of $p = \frac{m}{n}$ is much better, leading to a speed-up by a factor of order $m^{\Theta(m)}$.

Since [DLMN17] also observed that missing the optimal rate by a small constant factor leads to performance losses exponential in m , it was pro-

posed to use a mutation rate that is drawn from a (heavy-tailed) power-law distribution. Without the need to know m , this approach led to runtimes that exceed the ones obtained from the optimal rate $p = \frac{m}{n}$ by only a small factor polynomial in m . This price for universality can be made as low as $\Theta(m^{0.5+\varepsilon})$, but not smaller than $\Theta(\sqrt{m})$. Various variants of heavy-tailed mutation operators have been proposed subsequently, also heavy-tailed choices of other parameters have been used with great success [FQW18, FGQW18b, FGQW18a, WQT18, ABD20a, ABD20b, AD20, ABD21, DZ21, COY21].

A different way to cope with local optima was proposed in [RW20]. When an algorithm is stuck on a local optimum for a sufficiently long time, then with high probability it has explored all search points in a certain radius. Consequently, it is safe to increase the mutation rate, which increases the probability to generate more distant solutions. This is the main idea of a series of works on stagnation detection [RW20, RW21a, RW21b]. As shown in [RW20], this approach can save the polynomial price for universality of the heavy-tailed approach and thus obtain runtimes of the same asymptotic order as when using the optimal (problem-specific) mutation rate. By replacing standard-bit mutation with m -bit flips, the time to find a particular solution in Hamming distance m was further reduced to $(1 + o(1))\binom{n}{m}$, the same time (apart from lower order terms) one would obtain with the best unbiased mutation operator (which consists of flipping m random bits).

To be precise, two approaches are discussed in [RW21b]. The simple one, obtained from just replacing standard-bit mutation in [RW20] by r -bit mutation, obtains the desired runtimes with high probability, but fails completely with some very small probability. For this reason, also a robust version of the algorithm was proposed in [RW21b], which by cyclically reverting to smaller mutation strengths overcomes the problem that, with small probability, a given solution in distance m is not found in the phase which uses m -bit flips. In [RW21a], a variation of SD-RLS was proposed that keeps the successful strength after leaving local optima with the help of the radius memory mechanism, which is beneficial on highly multimodal fitness landscapes. The idea of stagnation detection has also been successfully used in multi-objective evolutionary computation [DZ21].

3 Algorithm SD-FEA $_{\beta,\gamma,R}$

We propose the algorithm SD-FEA $_{\beta,\gamma,R}$ for the maximization of pseudo-Boolean functions $f: \{0, 1\}^n \rightarrow \mathbb{R}$ defined in Algorithm 1. The function $\text{pow}(\beta, u)$ makes a sample from a power-law distribution with exponent β on $[1..u]$ as defined in 1 below.

Algorithm 1: SD-FEA $_{\beta,\gamma,R}$ for the maximization of $f: \{0, 1\}^n \rightarrow \mathbb{R}$

```

1 Select  $x$  uniformly at random from  $\{0, 1\}^n$  and set  $r_1 \leftarrow 1$ ;
2  $u \leftarrow 0$ ;
3 for  $t \leftarrow 1, 2, \dots$  do
    Set  $s = r_t$  with probability  $1 - \gamma$  or
4      $s = r_t + \text{pow}(\beta, n - r_t)$  with probability  $\gamma/2$  or
      $s = r_t - \text{pow}(\beta, \max\{1, r_t - 1\})$  with probability  $\gamma/2$ ;
5 Create  $y$  by flipping  $s$  bits in a copy of  $x$  uniformly;
6  $u \leftarrow u + 1$ ;
7 if  $f(y) > f(x)$  then
8      $x \leftarrow y$ ;
9      $r_{t+1} \leftarrow 1$ ;
10     $u \leftarrow 0$ ;
11 else if  $f(y) = f(x)$  and  $r_t = 1$  then
12      $x \leftarrow y$ ;
13 if  $u \geq \ell_{r_t}$  then
14      $r_{t+1} \leftarrow \min\{r_t + 1, \lfloor \frac{n}{2.1} \rfloor\}$ ;
15      $u \leftarrow 0$ ;
16 else
17      $r_{t+1} \leftarrow r_t$ ;

```

The general idea of this algorithm is that it increases the mutation strength r to $r + 1$ when the improvement is not at the Hamming distance smaller than r with at least a constant probability (with probability $1/R$ roughly) using the stagnation detection mechanism. Meanwhile, where the strength is r , called in phase r , the algorithm looks at larger or smaller Hamming distances (with probability γ) besides the current strength r . The distribution over the search distance to the current strength r follows a power-law distribution. An integer random variable X follows a power-law distribution with parameters β and u if

$$\Pr[X = i] = \begin{cases} C_{\beta,u} i^{-\beta} & \text{if } 1 \leq i \leq u, \\ 0 & \text{otherwise,} \end{cases} \quad (1)$$

where $C_{\beta,u} := (\sum_{j=1}^u j^{-\beta})^{-1}$ is the normalization coefficient. The function $\text{pow}(\beta, u)$ used in Algorithm 1 returns an integer from 1 to u sampled from the power-law distribution with parameters β and u .

The algorithm starts with a search point selected uniformly at random from the search space $\{0, 1\}^n$ with the initial strength $r = 1$. There is

a counter u for counting the number of unsuccessful steps with finding a strict improvement with the current strength. When the counter exceeds the threshold ℓ_r , the strength r increases by one if not beyond $n/2.1$, and when the algorithm makes progress, the counter and strength are reset to their initial values. The mutation, which we call s -flip in the following, flips exactly s bits randomly chosen as follows. With probability $1 - \gamma$, the algorithm flips exactly r bits in phase r . However, with probability γ , the algorithm flips less or more bits than r with the same chance with the help of power law distribution with parameter β . The distribution over s is analyzed in Lemma 1 below.

Regarding the parameters, in this paper, we use

$$\ell_r = \binom{n}{r} / (1 - \gamma) \ln(R). \quad (2)$$

This threshold almost fits for the pseudo-Boolean fitness functions. For other search spaces, the threshold should be $\ell_r = |S_r| / (1 - \gamma) \ln(R)$, where $|S_r|$ is the number of search points at the distance r from the current search point. The threshold defined in equation (2) has a parameter R controlling the probability of failing to find an improvement at the “right” strength. To prove the theoretical results, R should be selected at least $e^{1/\gamma}$. The parameter γ denotes the probability of having deviation for the number of bits flipped from the current strength r . In Section 7, we give some recommendations for choosing these parameters.

The *runtime* or the *optimization time* of a heuristic algorithm on a fitness function f is the first point of time t where a search point of maximal fitness has been selected.

4 Analysis of the Algorithm SD-FEA $_{\beta, \gamma, R}$

In this paper, we call the *fitness level gap* of a point $x \in \{0, 1\}^n$ the maximum of all individual gap sizes in the fitness level of x , where the *individual gap* is the minimum Hamming distance to points with strictly larger fitness function value, i. e.,

$$\begin{aligned} \text{IndividualGap}(x) &:= \min\{H(x, y) : f(y) > f(x), y \in \{0, 1\}^n\}, \\ \text{FitnessLevelGap}(x) &:= \max_{\{y | f(y) = f(x)\}} \text{IndividualGap}(y). \end{aligned}$$

If the algorithm creates a point at the Hamming distance $\text{FitnessLevelGap}(x)$ from the current search point x , with a positive probability, an improvement can be found. Note that $\text{FitnessLevelGap}(x) = 1$ is allowed, so the definition

also covers search points that are not local optima. As long as a strict improvement is not made, the FitnessLevelGap remains the same, although the current search point might be replaced with another search point in the fitness level in phase 1.

Let us define by the *epoch* of i the sequence of iterations with the search points in fitness level i . We define by phase r all points of time where radius r is used in the algorithm for an epoch. Let E_r be the event of **not** finding the optimum within phase r and for $j \geq i$ and E_i^j denote the event of not finding a strict improvement for the strengths i to j . Formally, $E_i^j = E_i \cap \dots \cap E_j$.

Before computing the probabilities of these events, we need to know the distribution of the offspring in an iteration. The following lemma will be used throughout this paper, showing the distribution of the number of flipping bits (i. e., the variable s in Algorithm 1) in each iteration. In phase r , with a relatively large probability $1 - \gamma$, the algorithm flips r bits. However, with probability γ , it uses power-law distributions to flip less or more than r bits.

Lemma 1. *Let r be the current strength in a run of the algorithm $SD\text{-}FEA_{\beta,\gamma,R}$. Let X be the integer random variable corresponding to the variable s in Algorithm 1, that is, the number of bits that are flipped. Then*

$$\Pr[X = \alpha] = \begin{cases} \gamma/2 \cdot C_{\beta,r-1} \cdot (r - \alpha)^{-\beta} & 1 \leq \alpha < r, \\ 1 - \gamma & \alpha = r, \\ \gamma/2 \cdot C_{\beta,n-r} \cdot (\alpha - r)^{-\beta} & \alpha > r. \end{cases}$$

Proof. It is immediately visible from Algorithm 1 that $\Pr[X = r] = 1 - \gamma$.

For $1 \leq \alpha < r$, we have

$$\begin{aligned} \Pr[X = \alpha] &= \Pr[X < r] \cdot \Pr[X = \alpha \mid X < r] \\ &= \Pr[X < r] \cdot \Pr[\text{pow}(\beta, r - 1) = r - \alpha] \\ &= \gamma/2 \cdot C_{\beta,r-1} (r - \alpha)^{-\beta}. \end{aligned}$$

For $\alpha > r$, we similarly obtain

$$\begin{aligned} \Pr[X = \alpha] &= \Pr[X > r] \cdot \Pr[X = \alpha \mid X > r] \\ &= \Pr[X > r] \cdot \Pr[\text{pow}(\beta, n - r) = \alpha - r] \\ &= \gamma/2 \cdot C_{\beta,n-r} (\alpha - r)^{-\beta}. \end{aligned}$$

□

The following lemma shows that the probability of reaching a phase r is greater than the fitness gap size. In the statement of the lemma, recall that the parameter R controls the length of the phase.

Lemma 2. Let $x \in \{0, 1\}^n$ be the current search point of $SD\text{-}FEA_{\beta, \gamma, R}$ with $\beta > 1$ on a pseudo-Boolean fitness function $f: \{0, 1\}^n \rightarrow \mathbb{R}$, and $m = \text{FitnessLevelGap}(x)$. Then for $m < r \leq \lfloor \frac{n}{2.1} \rfloor$, we have

$$\Pr[E_1^{r-1}] \leq R^{-1-\gamma/2 \cdot \left(\frac{\ln(1.1)}{\beta}\right)^\beta C_{\beta, n}(r-m-1)},$$

where E_1^{r-1} denotes the probability of not finding an improvement in phases 1 to $r-1$.

Proof. Let p_r be a lower bound on the probability of making progress in phase r in one iteration. Then we have

$$\begin{aligned} \Pr[E_1^{r-1}] &\leq \Pr[E_m \cap \dots \cap E_{r-1}] = \prod_{i=m}^{r-1} \Pr[E_i] \leq \prod_{i=m}^{r-1} (1-p_i)^{\ell_i} \\ &\leq \exp\left(-\sum_{i=m}^{r-1} p_i \ell_i\right), \end{aligned} \quad (3)$$

where we use the inequality $1+x \leq e^x$ for all $x \in \mathbb{R}$.

In the following paragraphs, we aim at bounding $p_i \ell_i$ from below.

For $i = m$, via Lemma 1 and since $\ell_r = \binom{n}{r}/(1-\gamma)\ln(R)$, we have

$$p_m \ell_m \geq (1-\gamma) \binom{n}{m}^{-1} \cdot \binom{n}{m} / (1-\gamma) \ln(R) = \ln(R).$$

For $m < i \leq \frac{n}{2.1}$, again using Lemma 1, we have

$$p_i \geq \gamma/2 C_{\beta, i-1} (i-m)^{-\beta} \binom{n}{m}^{-1}.$$

Then we bound $p_i \ell_i$ from below by

$$p_i \ell_i \geq \gamma/2 \cdot C_{\beta, i-1} \frac{\binom{n}{i}/(1-\gamma)\ln(R)}{(i-m)^\beta \binom{n}{m}} \geq \gamma/2 \cdot C_{\beta, n} \frac{\binom{n}{i} \ln(R)}{(i-m)^\beta \binom{n}{m}},$$

where we have $C_{\beta, n} \leq C_{\beta, i-1}$. The last expression is bounded from below by

$$\gamma/2 \cdot C_{\beta, n} \frac{\ln(R)}{(i-m)^\beta} \cdot \frac{\binom{n}{i}}{\binom{n}{i-1}} \dots \frac{\binom{n}{m+1}}{\binom{n}{m}} \geq \gamma/2 \cdot C_{\beta, n} \frac{(1.1)^{i-m} \ln(R)}{(i-m)^\beta},$$

where we have $\binom{n}{k}/\binom{n}{k-1} = \frac{n-k+1}{k} \geq 1.1$ for $k \leq \lfloor \frac{n}{2.1} \rfloor$. Also, we claim that $1.1^k/k^\beta \geq (\ln(1.1)/\beta)^\beta$ for $k \in \mathbb{N}_{\geq 1}$. To prove, let $f(x) = 1.1^x/x^\beta$. For $x > 0$,

its derivative, i. e., $f'(x)$, has only one root $\hat{x} = \frac{\beta}{\ln 1.1}$. Before and after this point the function is decreasing and increasing, respectively, so $f(\hat{x})$ is the minimum value of the function for $x > 0$. We have

$$f(\hat{x}) = \frac{1.1^{\beta/\ln(1.1)}}{(\beta/\ln(1.1))^\beta} \geq \left(\frac{\ln(1.1)}{\beta}\right)^\beta.$$

Thus, the last expression is bounded from below by $\gamma/2 \cdot C_{\beta,n}(\ln(1.1)/\beta)^\beta \ln(R)$.

Finally, from Equation (3), we obtain

$$\Pr[E_1^{r-1}] \leq \exp\left(-\sum_{i=m}^{r-1} p_i \ell_i\right) \leq R^{-1-\gamma/2 \cdot \left(\frac{\ln(1.1)}{\beta}\right)^\beta C_{\beta,n}(r-m-1)}.$$

□

The next lemma is used to estimate the number of iterations in phases larger than the fitness level gap. With a good choice of the parameters γ and R , the following result can be improved to $o(1/s_m)$, that is to say the number of steps at larger strengths can be captured by the number of steps at the phase m .

Lemma 3. *Let $x \in \{0, 1\}^n$ be the current search point of $SD\text{-}FEA_{\beta,\gamma,R}$ with $\beta > 1$ and $R \geq e^{1/\gamma}$ on a pseudo-Boolean function $f: \{0, 1\}^n \rightarrow \mathbb{R}$. Assume $m = \text{FitnessLevelGap}(x)$ and $m \leq \lfloor n/2.1 \rfloor$. Let s_m be a lower bound on the probability that an improvement is found from search points in the fitness level of x conditional on flipping m bits. Then, the expected number of iterations spent in larger strengths than m , i. e., $E[I_{>m}]$, is at most*

$$O\left(R^{-1}\gamma^{-1}\frac{1}{s_m}\right).$$

Proof. Let I_r be the number of iterations spent in phase r . Then

$$E[I_{>m}] = \sum_{r=m+1}^{\lfloor \frac{n}{2.1} \rfloor - 1} E[I_r] + E[I_{\lfloor \frac{n}{2.1} \rfloor}].$$

With probability $\Pr[E_1^{r-1}]$, the algorithm does not make progress with strengths less than r . In phase r , the probability of finding an improvement is $C_{\beta,r-1}\gamma/2(r-m)^{-\beta} \cdot s_m$ in each iteration using Lemma 1. Thus, for all strengths $r > m$, using the law of total probability, we have

$$\begin{aligned} E[I_r] &= \Pr[E_1^{r-1}] E[I_r | E_1^{r-1}] + \Pr[\overline{E_1^{r-1}}] E[I_r | \overline{E_1^{r-1}}] \\ &\leq \Pr[E_1^{r-1}] \cdot (C_{\beta,r-1})^{-1} 2/\gamma \cdot \frac{1}{s_m} (r-m)^\beta + 0. \end{aligned}$$

Using Lemma 2 and since $R \geq e^{1/\gamma}$, we can bound

$$\begin{aligned} E[I_r] &\leq R^{-1-\gamma/2 \cdot (\frac{\ln(1.1)}{\beta})^\beta} C_{\beta,n}(r-m-1) 2/\gamma \cdot \frac{1}{s_m} (r-m)^\beta \\ &= O\left(R^{-1}\gamma^{-1} \frac{1}{s_m} \frac{(r-m)^\beta}{e^{1/2 \cdot (\frac{\ln(1.1)}{\beta})^\beta} C_{\beta,n}(r-m-1)} \right), \end{aligned}$$

where we have $(C_{\beta,r-1})^{-1} = O(1)$ for $\beta > 1$. This results in

$$\begin{aligned} \sum_{r=m+1}^{\lfloor \frac{n}{2.1} \rfloor - 1} E[I_r] &\leq O\left(R^{-1}\gamma^{-1} \frac{1}{s_m} \sum_{r=m+1}^{\lfloor \frac{n}{2.1} \rfloor - 1} \frac{(r-m)^\beta}{e^{1/2 \cdot (\frac{\ln(1.1)}{\beta})^\beta} C_{\beta,n}(r-m)} \right) \\ &\leq O\left(R^{-1}\gamma^{-1} \frac{1}{s_m} \right), \end{aligned}$$

where we use

$$\sum_{r=m+1}^{\lfloor \frac{n}{2.1} \rfloor - 1} \frac{(r-m)^\beta}{e^{1/2 \cdot (\frac{\ln(1.1)}{\beta})^\beta} C_{\beta,n}(r-m)} = \sum_{r=m+1}^{\lfloor \frac{n}{2.1} \rfloor - 1} \frac{(r-m)^\beta}{e^{\Theta(r-m)}} = O(1).$$

Similarly, we have

$$E[I_{\lfloor \frac{n}{2.1} \rfloor}] \leq \Pr[E_1^{\lfloor \frac{n}{2.1} \rfloor - 1}] \cdot \gamma/2 \cdot \frac{1}{s_m} (\lfloor \frac{n}{2.1} \rfloor - m)^\beta \leq O\left(R^{-1}\gamma^{-1} \frac{1}{s_m} \right).$$

Altogether, we obtain

$$E[I_{>m}] = \sum_{r=m+1}^{\lfloor \frac{n}{2.1} \rfloor - 1} E[I_r] + E[I_{\lfloor \frac{n}{2.1} \rfloor}] = O\left(R^{-1}\gamma^{-1} \frac{1}{s_m} \right),$$

as suggested. □

In the following lemma, which has been taken from [RW21b], we have a combinatorial inequality that will be used in the analyses of the algorithms to count the number of iterations spent in smaller strengths than the fitness level gap.

Lemma 4 (Lemma 1 in [RW21b]). *For any integer $m \leq n/2$, we have*

$$\sum_{i=1}^m \binom{n}{i} \leq \frac{n - (m-1)}{n - (2m-1)} \binom{n}{m}.$$

We now present the first main result. In the following theorem, given $m \leq n/2.1$, we provide a rigorous upper bound on the escaping time from a local optimum. We also prove a general bound on the expected time to leave it for an arbitrary gap size.

Theorem 5. *Let $x \in \{0, 1\}^n$ be the current search point. Define T as the time $SD\text{-}FEA_{\beta, \gamma, R}$ with $\beta > 1$ and $R \geq e^{1/\gamma}$ takes to create a strict improvement on a pseudo-Boolean function $f: \{0, 1\}^n \rightarrow \mathbb{R}$. Assume $m = \text{FitnessLevelGap}(x)$. If $m \leq n/2.1$, we have*

$$E[T] \leq \binom{n}{m} \left(\frac{1}{1-\gamma} + O\left(\frac{m \ln(R)}{(1-\gamma)n} + R^{-1}\gamma^{-1}\right) \right).$$

Moreover, for all $m \leq n$, we have

$$E[T] = O\left(2^n \ln(R) + 1/\gamma \binom{n}{m} \lfloor \frac{n}{2.1} \rfloor - m \right)^\beta.$$

Proof. Let I_r be the number of iterations spent in phase r . Using linearity of expectation, we have

$$E[T] = \sum_{r=1}^{\lfloor \frac{n}{2.1} \rfloor - 1} E[I_r] + E[I_{\lfloor \frac{n}{2.1} \rfloor}].$$

Regarding the first part, we have $m \leq n/2.1$. If the phase is less than m , i. e., $r < m$, we have that $E[I_r]$ is at most the threshold value at phase r , i. e., $\ell_r = \binom{n}{r} / (1-\gamma) \ln(R)$. Thus, by using Lemma 4, we compute

$$\begin{aligned} \sum_{r=1}^{m-1} E[I_r] &\leq \sum_{r=1}^{m-1} \binom{n}{r} \frac{\ln(R)}{1-\gamma} \\ &\leq \binom{n}{m-1} \frac{\ln(R)}{1-\gamma} \cdot \frac{n-(m-2)}{n-(2m-3)} \\ &= \binom{n}{m} \frac{\ln(R)}{1-\gamma} \cdot \frac{m}{n-m+1} \cdot \frac{n-(m-2)}{n-(2m-3)}. \end{aligned}$$

Since $m \leq \frac{n}{2.1}$, the last expression is bounded from above by

$$\sum_{r=1}^{m-1} E[I_r] = O\left(\binom{n}{m} \frac{m \ln(R)}{(1-\gamma)n}\right).$$

When the strength is m , with probability $1-\gamma$, the algorithm flips exactly m bits (Lemma 1). When m bits are flipped, with probability $\binom{n}{m}^{-1}$,

an improvement is found. Regarding a truncated geometric distribution with success probability $(1-\gamma)\binom{n}{m}^{-1}$, within $\binom{n}{m}/(1-\gamma)$ iterations in expectation, we see that the algorithm finds a better point or the phase is terminated. Thus,

$$E[I_m] \leq \frac{\binom{n}{m}}{(1-\gamma)}.$$

For $r > m$, using Lemma 3 with $s_m \geq \binom{n}{m}^{-1}$, we obtain

$$E[I_{>m}] = \sum_{r=m+1}^{\lfloor \frac{n}{2.1} \rfloor - 1} E[I_r] + E[I_{\lfloor \frac{n}{2.1} \rfloor}] = O\left(R^{-1}\gamma^{-1}\binom{n}{m}\right).$$

Altogether, we have

$$\begin{aligned} E[T] &= \sum_{r=1}^{\lfloor \frac{n}{2.1} \rfloor - 1} E[I_r] + E[I_{\lfloor \frac{n}{2.1} \rfloor}] \\ &= \sum_{r=1}^{m-1} E[I_r] + E[I_m] + \sum_{r=m+1}^{\lfloor \frac{n}{2.1} \rfloor - 1} E[I_r] + E[I_{\lfloor \frac{n}{2.1} \rfloor}] \\ &\leq \binom{n}{m} \left(\frac{1}{1-\gamma} + O\left(\frac{m \ln(R)}{(1-\gamma)n} + R^{-1}\gamma^{-1}\right) \right). \end{aligned}$$

Regarding the second part, since for $r \leq \lfloor \frac{n}{2.1} \rfloor - 1$, we have that $E[I_r]$ is at most the threshold value, we have

$$E[T] \leq \sum_{r=1}^{\lfloor \frac{n}{2.1} \rfloor - 1} \ell_r + E[I_{\lfloor \frac{n}{2.1} \rfloor}] = \sum_{r=1}^{\lfloor \frac{n}{2.1} \rfloor - 1} \binom{n}{r} / (1-\gamma) \ln(R) + E[I_{\lfloor \frac{n}{2.1} \rfloor}].$$

In phase $\lfloor \frac{n}{2.1} \rfloor$, the algorithm no longer increases the strength until finding an improvement. Using Lemma 1, the improvement is found with probability at least

$$\Omega\left(\gamma/2 \cdot \left|\lfloor \frac{n}{2.1} \rfloor - m\right|^{-\beta} \cdot \binom{n}{m}^{-1}\right)$$

in each iteration. Using the geometric distribution with this success probability, we obtain

$$\begin{aligned} E[T] &\leq \sum_{r=1}^{\lfloor \frac{n}{2.1} \rfloor - 1} \binom{n}{r} / (1-\gamma) \ln(R) + O\left(1/\gamma \binom{n}{m} \left|\lfloor \frac{n}{2.1} \rfloor - m\right|^\beta\right) \\ &= O\left(2^n \frac{\ln(R)}{1-\gamma} + 1/\gamma \binom{n}{m} \left|\lfloor \frac{n}{2.1} \rfloor - m\right|^\beta\right), \end{aligned}$$

where we have $\sum_{i=1}^n \binom{n}{i} = 2^n$. The second part is proved as desired. \square

Theorem 5 provides a precise upper bound on the escaping time from local optima where there are a few ways to leave them. However, it is not practical when there are significantly more ways to jump over the valley from the local optima. The following theorem considers such scenarios. The constant r' defined in the theorem basically represents the first phase that the probability of finding one of the improvements is at least constant, and its value is an integer between 1 and δ .

Theorem 6. *Let $x \in \{0, 1\}^n$ be the current search point and s_m be a lower bound on the probability that a strict improvement is found from search points in the fitness level of x conditional on flipping m bits. Define T as the time $SD\text{-}FEA_{\beta, \gamma, R}$ with $\beta > 1$ and $R \geq e^{1/\gamma}$ takes to create a strict improvement on a pseudo-Boolean function $f: \{0, 1\}^n \rightarrow \mathbb{R}$. Assume $m = \text{FitnessLevelGap}(x)$. Then for $m \leq n/2.1$, we have*

$$E[T] \leq \frac{1}{s_m} \cdot \frac{1}{\gamma} (\delta - r')^\beta \cdot O\left(1 + \frac{r' \ln(R)}{(1 - \gamma)n} + R^{-1}\right).$$

where $r' = \min\left\{\delta, \arg \max_r \left\{\binom{n}{r} \leq \frac{1}{s_m} \frac{1}{\gamma} (\delta - r)^\beta\right\}\right\}$.

Proof. Let I_r be the number of iterations spent in phase r . using linearity of expectation, we have

$$E[T] = \sum_{r=1}^{\lfloor \frac{n}{2.1} \rfloor - 1} E[I_r] + E[I_{\lfloor \frac{n}{2.1} \rfloor}].$$

If the strength is less than r' , i. e., $r < r'$, we have that $E[I_r]$ is at most the threshold value at phase r . Thus, we have

$$\begin{aligned} \sum_{r=1}^{r'-1} E[I_r] &\leq \sum_{r=1}^{r'-1} \binom{n}{r} / (1 - \gamma) \ln(R) \\ &\leq \binom{n}{r'-1} \frac{\ln(R)}{(1 - \gamma)} \frac{n - (r' - 2)}{n - (2r' - 3)} \\ &= \frac{r'}{n - r' + 1} \cdot \binom{n}{r'} \frac{\ln(R)}{(1 - \gamma)} \frac{n - (r' - 2)}{n - (2r' - 3)}. \end{aligned}$$

Since $r' \leq m \leq \frac{n}{2.1}$, the last expression is bounded from above by

$$\sum_{r=1}^{r'-1} E[I_r] = O\left(\binom{n}{r'} \frac{r' \ln(R)}{(1 - \gamma)n}\right).$$

Using the definition of r' described in the theorem statement, we get

$$\sum_{r=1}^{r'-1} E[I_r] = O\left(\frac{1}{s_m} \cdot \frac{1}{\gamma} (\delta - r')^\beta \cdot \frac{r' \ln(R)}{(1-\gamma)n}\right).$$

In the phases from r' to $m-1$, the probability of finding an improvement is at least $s_m \cdot \gamma / 2C_{\beta, n-r'} (\delta - r')^{-\beta}$ (Lemma 1). Using the geometric distribution with this success probability, the success happens within

$$O\left(\frac{1}{s_m} \cdot \frac{1}{\gamma} (\delta - r')^\beta\right)$$

iterations in expectation.

In phase m , where the strength is m , exactly m bits are flipped with probability $1-\gamma$ (Lemma 1), and an improvement is found with probability at least s when m bits are flipped. Regarding a truncated geometric distribution with success probability $(1-\gamma)s_m$, within $1/s_m \cdot 1/(1-\gamma)$ iterations in expectation, we observe that the algorithm finds an improvement or the phase is terminated. Thus,

$$E[I_m] \leq \frac{1}{s_m} \cdot \frac{1}{(1-\gamma)}.$$

For $r > m$, using Lemma 3 with s_m , we obtain

$$E[I_{>m}] = O(R^{-1} \gamma^{-1} s_m^{-1}).$$

Altogether, we have

$$\begin{aligned} E[T] &= \sum_{r=1}^{\lfloor \frac{n}{2.1} \rfloor - 1} E[I_r] + E[I_{\lfloor \frac{n}{2.1} \rfloor}] \\ &= \sum_{r=1}^{r'-1} E[I_r] + \sum_{r=r'}^{m-1} E[I_r] + E[I_m] + \sum_{r=m+1}^{\lfloor \frac{n}{2.1} \rfloor - 1} E[I_r] + E[I_{\lfloor \frac{n}{2.1} \rfloor}] \\ &\leq O\left(\frac{1}{s_m} \cdot \frac{1}{\gamma} (\delta - r')^\beta \cdot \frac{r' \ln(R)}{(1-\gamma)n} + \frac{1}{s_m} \cdot \frac{1}{\gamma} (\delta - r')^\beta + \frac{1}{s_m(1-\gamma)} + \frac{R^{-1}}{\gamma s_m}\right) \\ &\leq \frac{1}{s_m} \cdot \frac{1}{\gamma} (\delta - r')^\beta \cdot O\left(1 + \frac{r' \ln(R)}{(1-\gamma)n} + R^{-1}\right). \end{aligned}$$

□

After establishing some tools for obtaining upper bounds on the time required to escape from local optima, we aim to demonstrate the performance of $\text{SD-FEA}_{\beta,\gamma,R}$ on the sub-problems without local optima. On unimodal functions, the gap of all search points in the search space (except for the global optima) is 1, so the algorithm can make progress in phase 1. The following benchmark functions ONEMAX and LEADINGONES have been extensively studied in the literature as unimodal problems.

$$\text{ONEMAX}(x_1, \dots, x_n) := \|x\|_1 \text{ and } \text{LEADINGONES}(x_1, \dots, x_n) := \sum_{i=1}^n \prod_{j=1}^i x_j,$$

where $\|x\|_1$ is the number of one-bits in the bit string.

In the following theorem, we state how $\text{SD-FEA}_{\beta,\gamma,R}$ behaves on unimodal functions compared to RLS using an upper bound based on the fitness-level method [Weg01].

Theorem 7. *Let $f: \{0, 1\}^n \rightarrow \mathbb{R}$ be a unimodal function and $|\text{Im}(f)|$ be the number of fitness values of the underlying function f . Let f_i be the i -th fitness value of an increasing order of all fitness values in f . We consider all fitness levels $A_1, \dots, A_{|\text{Im}(f)|}$ such that A_i contains search points with fitness value f_i . Let s_i be a lower bound on the probability that RLS finds an improvement from the search points in fitness level A_i . Denote by T the runtime of $\text{SD-FEA}_{\beta,\gamma,R}$ with $\beta > 1$ and $R \geq e^{1/\gamma}$ on f . Then*

$$E[T] \leq \left(\frac{1}{1-\gamma} + O(R^{-1}\gamma^{-1}) \right) \sum_{i=1}^{|\text{Im}(f)|-1} \frac{1}{s_i}.$$

Proof. We define by $I^{(i)}$ the number of all iterations spent to leave the fitness level i . using linearity of expectation, we have

$$E[T] = \sum_{i=1}^{|\text{Im}(f)|-1} E[I^{(i)}].$$

Let $I_r^{(i)}$ be the number of iterations spent in phase r after a search point for A_i was found. Then we define

$$I^{(i)} = \sum_{r=1}^{\lfloor \frac{n}{2-1} \rfloor - 1} I_r^{(i)} + I_{\lfloor \frac{n}{2-1} \rfloor}^{(i)}.$$

As long as the strength is 1, the algorithm flips exactly one-bit with probability at least $1 - \gamma$ (Lemma 1). The worst-case time to leave fitness

level i is at most $1/(1-\gamma) \cdot 1/s_i$ using the geometric distribution with success probability $s_i \cdot (1-\gamma)$. Hence, for each fitness level i , we bound $I_1^{(i)}$ from above by $1/(1-\gamma) \cdot 1/s_i$, and for $r > 1$, we bound $I_r^{(i)}$ from above by using Lemma 3 with $s_m = s_i$. Hence,

$$E[I_{>1}^{(i)}] = O\left(R^{-1}\gamma^{-1}\frac{1}{s_i}\right),$$

Altogether, we have

$$\begin{aligned} E[T] &= \sum_{i=1}^{|\text{Im}(f)|-1} E[I^{(i)}] \\ &\leq \sum_{i=1}^{|\text{Im}(f)|-1} \left(\frac{1}{s_i(1-\gamma)} + O\left(R^{-1}\gamma^{-1}\frac{1}{s_i}\right) \right) \\ &\leq \left(\frac{1}{1-\gamma} + O(R^{-1}\gamma^{-1}) \right) \sum_{i=1}^{|\text{Im}(f)|-1} \frac{1}{s_i}. \end{aligned}$$

□

The corollary below is a result of Theorem 7 applied on the unimodal functions ONEMAX with $s_i = (n - (i - 1))/n$ and LEADINGONES with $s_i = 1/n$.

Corollary 8. *The expected runtime of the SD-FEA $_{\beta,\gamma,R}$ with $\beta > 1$, $\gamma = o(1)$ and $R \geq \Omega(e^{1/\gamma})$ on ONEMAX is at most $(1 + o(1))n \ln n$ and on LEADINGONES is at most $(1 + o(1))n^2$.*

5 Analysis on Jump $_{k,\delta}$

In this section, we use the results in the previous section to prove a bound on a generalization of JUMP $_{\delta}$ called JUMP $_{k,\delta}$ with two parameters k and δ , see Figure 1 for a depiction.

This function is based on the well-known JUMP benchmark [DJW02], in which the place of the jump with size δ starts at the Hamming distance k from the global optimum. In other words, after the jump, there is a unimodal sub-problem of length $k - \delta$. The classical JUMP function is a special case of JUMP $_{k,\delta}$ with $k = \delta$, i. e., JUMP $_{\delta} = \text{JUMP}_{\delta,\delta}$. Formally,

$$\text{JUMP}_{k,\delta}(x) = \begin{cases} \|x\|_1 & \text{if } \|x\|_1 \in [0..n - k] \cup [n - k + \delta..n], \\ -\|x\|_1 & \text{otherwise.} \end{cases}$$

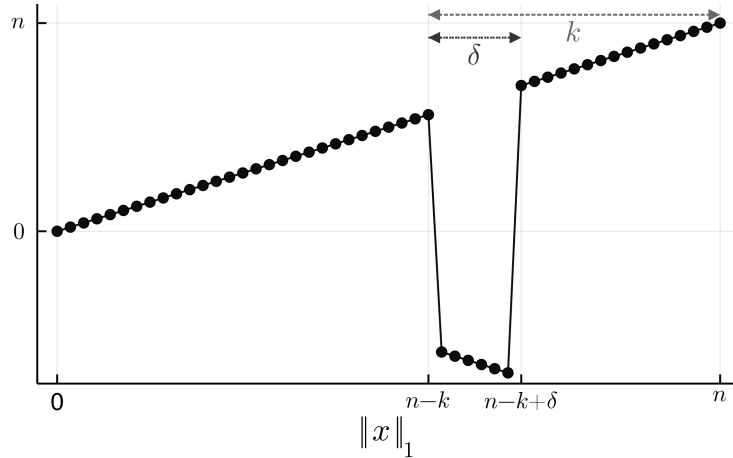


Figure 1: The function $\text{JUMP}_{k,\delta}$.

We refer the interested reader to see [BBD21] for more information about $\text{JUMP}_{k,\delta}$, where the performance of the $(1+1)$ EA, the $(1+1)$ FEA_β , and the robust version of SD-RLS (SD-RLS^r) are carefully analyzed. Also, Rajabi and Witt [RW21a] independently define the jump function with an offset to analyze the recovery time for the strength in the algorithm SD-RLS with radius memory (SD-RLS^m) after leaving the local optimum. Recently, Witt in [Wit21] analyzes the performance of some other algorithms on the function $\text{JUMP}_{k,\delta}$ (which is so-called JUMPOFFSET in the paper).

We want to show that the algorithm $\text{SD-FEA}_{\beta,\gamma,R}$ performs relatively efficiently on $\text{JUMP}_{k,\delta}$ in both cases when $k = \delta$ (i.e., JUMP_δ) and $k > \delta$. In the first case, when basically, there are not many improving solutions, $\text{SD-FEA}_{\beta,\gamma,R}$ with $\gamma = o(1)$ optimizes JUMP_δ as efficient as SD-RLS^r thanks to Theorem 5. The result is formally proved in Theorem 9.

Theorem 9. *The expected runtime $E[T]$ of $\text{SD-FEA}_{\beta,\gamma,R}$ with $\beta > 1$, $\gamma = o(1)$ and $R = \Omega(e^{1/\gamma})$ on JUMP_δ with $2 \leq \delta = o(n/\ln(R))$ satisfies*

$$E[T] \leq \binom{n}{\delta} (1 + o(1)).$$

Proof. Before reaching a local optimum with $n - m$ one-bits, JUMP_δ is equivalent to ONEMAX . Thus, the expected time until $\text{SD-FEA}_{\beta,\gamma,R}$ reaches the local optimum is at most $O(n \ln n)$ via Theorem 7 with $s_i = (n - (i - 1))/n$.

For a local optimum x we have $\text{FitnessLevelGap}(x) = \delta$ according to the definition of JUMP . Hence, using Theorem 5, the algorithm finds the global

optimum from the local optimum within the expected time at most

$$\binom{n}{\delta}(1 + o(1)).$$

This dominates the expected time of the algorithm before the local optimum. \square

It is also easy to see (similarly to the analysis of Theorem 9) that for $\gamma = \Theta(1)$, the expected runtime of SD-FEA $_{\beta,\gamma,R}$ is

$$\binom{n}{\delta} \left(\frac{1}{1-\gamma} + o(1) \right),$$

which is still asymptotically efficient.

We now present an upper bound on the optimization time of the proposed algorithm on JUMP $_{k,\delta}$.

Theorem 10. *The expected runtime $E[T]$ of SD-FEA $_{\beta,\gamma,R}$ with $\beta > 1$, $R = \Omega(e^\gamma)$ on JUMP $_{k,\delta}$ with $\delta = o(n/\ln(R))$ satisfies*

$$E[T] = O \left(\binom{n}{\delta} \binom{k}{\delta}^{-1} (\delta - r')^\beta \cdot \gamma^{-1} + n \ln n \right),$$

where $r' = \min \left\{ \delta, \arg \max_r \left\{ \binom{n}{r} \leq \binom{n}{\delta} \binom{k}{\delta}^{-1} \frac{1}{\gamma} (\delta - r)^\beta \right\} \right\}$.

Proof. Until reaching the local optimum with $n - k$ one-bits, JUMP $_{k,\delta}$ is equivalent to ONEMAX. Thus, the expected time until SD-FEA $_{\beta,\gamma,R}$ reaches the local optimum is at most $O(n \ln n)$ via Theorem 7 with $s_i = (n - (i - 1))/n$.

For a local optimum x , we have $\text{FitnessLevelGap}(x) = \delta$ according to the definition of JUMP $_{k,\delta}$. Using Theorem 5 with $s = \binom{n}{\delta}^{-1} \binom{k}{\delta}$, the algorithm finds a strict improvement with at least $n - k + \delta$ one-bits from the local optimum within expected time at most

$$O \left(\binom{n}{\delta} \binom{k}{\delta}^{-1} (\delta - r')^\beta \cdot \gamma^{-1} \right).$$

After leaving the local optimum, JUMP $_{k,\delta}$ is again equivalent to ONEMAX on the second slope. Using the same arguments as in the beginning of the proof, the expected time until SD-FEA $_{\beta,\gamma,R}$ reaches the global optimum is at most $O(n \ln n)$ via Theorem 7 with $s_i = (n - (i - 1))/n$. \square

In the following corollary, we see a scenario where we have $r' \geq \delta - c$ for some constant c , resulting in the term $(\delta - r')^\beta$ disappearing from the asymptotic upper bound. This is also an example where the SD-FEA $_{\beta,\gamma,R}$ can asymptotically outperform the (1+1) FEA $_\beta$.

Corollary 11. *Let $\Delta \geq 2$ be a constant. The expected runtime $E[T]$ of SD-FEA $_{\beta,\gamma,R}$ with $\beta > 1$ and $R = \Omega(e^\gamma)$ on JUMP $_{k,\delta}$ with $\omega(1) = k \leq \ln n$ and $\delta = k - \Delta$ satisfies*

$$E[T] = O\left(\binom{n}{\delta} \binom{k}{\delta}^{-1} \cdot \gamma^{-1}\right).$$

Proof. We claim that r' defined in Theorem 10 is at least $k - 2\Delta$. We compute

$$\begin{aligned} \frac{\binom{n}{k-2\Delta}}{\binom{n}{k-\Delta} \gamma^{-1} \Delta^\beta} &\leq \frac{(en/(k-2\Delta))^{k-2\Delta}}{n^{k-\Delta} k^{-\Delta} (k-\Delta)^{\Delta-k} \cdot \gamma^{-1} \Delta^\beta} \\ &= \gamma \cdot \frac{e^{k-2\Delta} \cdot k^\Delta \cdot (k-\Delta)^\Delta}{n^\Delta \Delta^\beta} \cdot \left(1 + \frac{\Delta}{k-2\Delta}\right)^{k-2\Delta} \\ &\leq \gamma \cdot \frac{e^{k-\Delta} \cdot k^\Delta \cdot (k-\Delta)^\Delta}{n^\Delta \Delta^\beta} = o(1), \end{aligned}$$

where we use the assumption $k \leq \ln n$, the inequality $(n/m)^m \leq \binom{n}{m} \leq (en/m)^m$ and $1 + x \leq e^x$ for all $x \in \mathbb{R}$. For a large enough n , the last expression results in

$$\binom{n}{k-2\Delta} \leq \binom{n}{k-\Delta} \binom{k}{\Delta}^{-1} \gamma^{-1} \Delta^\beta,$$

which means that $r' \geq k - 2\Delta$. Therefore, using the result of Theorem 10 with $r' \geq k - 2\Delta$, we get

$$E[T] = O\left(\binom{n}{\delta} \binom{k}{\delta}^{-1} \Delta^\beta \cdot \gamma^{-1} + n \ln n\right) = O\left(\binom{n}{\delta} \binom{k}{\delta}^{-1} \gamma^{-1}\right),$$

where $O(n \ln n)$ is captured by the first term according to our assumptions. \square

6 Experiments

In this section, we present the results of the experiments carried out to measure the performance of the proposed algorithm and several related ones on concrete problem sizes.

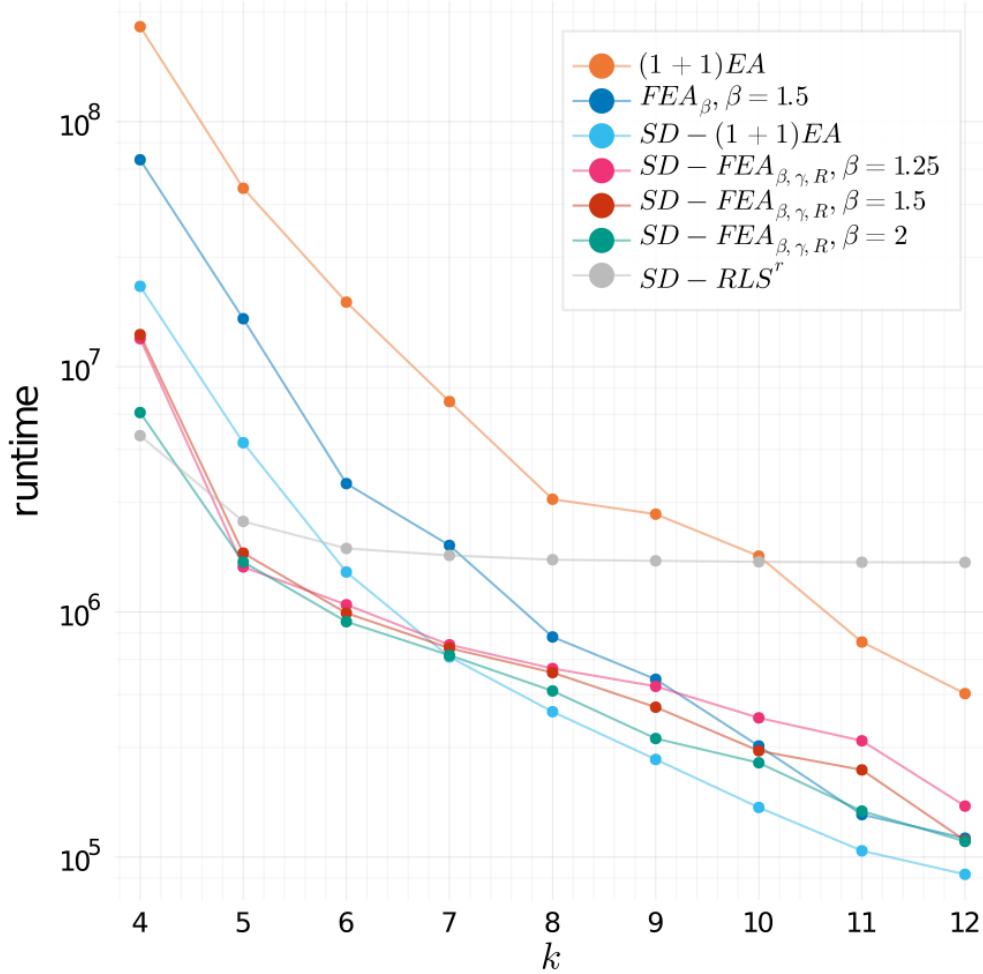


Figure 2: Average number (over 200 runs) of fitness calls the mentioned algorithms spent to optimize $JUMP_{k,4}$ with different values for k .

We ran an implementation of $SD-FEA_{\beta, \gamma, R}$ with $\beta \in \{1.25, 1.5, 2\}$, $\gamma = 1/4$ and $R = 25$ on the fitness function $JUMP_{k, \delta}$ of size $n = 100$ with the jump size $\delta = 4$ and k varying from 4 to 13. We recall that we have the classical JUMP function by setting $k = 4$. We compared our algorithm with the classical $(1+1)$ EA with standard mutation rate $1/n$, the $(1+1)$ FEA_{β} from [DLMN17] with $\beta = 1.5$, the $SD-(1+1)$ EA presented in [RW20] with $R = n^2$, and $SD-RLS^r$ from [RW21b] with $R = n^2$. The data presented in the figures is the average number of fitness calls over 200 runs. Considering the hypothesis of identical behavior, we use the Mann-

Whitney U-test between the algorithms, with the result that all p-values are less than 10^{-3} .

As can be seen in Figure 2, SD-RLS^r outperforms the rest of the algorithms where $k = 4$, i. e., there is only one improving solution for local optima. Our SD-FEA _{β, γ, R} needs around $(1 - \gamma)^{-1}$ times more fitness function calls than that, but still significantly outperforms the (1+1) FEA _{β} , SD-(1+1) EA and (1 + 1) EA. As k is increasing, the average running time of SD-RLS^r improves little and remains almost without change after $k = 5$; consequently, this algorithm becomes less and less competitive for growing k . This is natural since this algorithm necessarily has to reach phase 4 to be able to flip 4 bits. All other algorithms, especially the (1+1) FEA _{β} , perform increasingly better with larger k .

In a middle regime of $k \in \{5, 6, 7\}$, the SD-FEA _{β, γ, R} has the best average running time among the algorithms regarded. Although both with $k = 4$ and for $k \geq 8$, the SD-FEA _{β, γ, R} is not the absolutely best algorithm, but its performance loss over the most efficient algorithm (SD-RLS^r for $k = 4$ and SD-(1+1) EA for $k \geq 8$) is small. This finding supports our claim that our algorithm is a good approach to leaving local optima of various kinds.

For a large k , such as 10 or 11, the good performance of the SD-(1+1) EA and (1+1) FEA _{β} might appear surprising. The reason for the slightly weaker performance of our algorithm is the relatively small width of the valley of low fitness ($\delta = 4$), where our algorithm cannot fully show its advantages, but pays the price of sampling from the right heavy-tailed distribution only with probability $\frac{\gamma}{2}$.

7 Recommended Parameters

In this section, we use our theoretical and experimental results to derive some recommendations for choosing the parameters β , γ , and R of our algorithm. We note that having three parameters for a simple (1 + 1)-type optimizer might look frightening at first, but a closer look reveals that setting these parameters is actually not too critical.

For the power-law exponent β , as in [DLMN17], there is no indication that using a value different from $\beta = 1.5$ can give significant performance gains (in our experiments, $\beta = 2$ gave slightly better results, but we suspect that for larger jump sizes – not done here for reasons of computation time – this advantage vanishes). So clearly, this is the last parameter to try to optimize.

Different from the previous approaches building on stagnation detection, our algorithm also does not need specific values for the parameter R , which

governs the phase length $\ell_r = \frac{1}{1-\gamma} \binom{n}{r} \ln(R)$ and in particular leads to the property that a single improving solution in distance m is found in phase m with probability $1 - \frac{1}{R}$ (as follows from the proof of Lemma 2). Since we have the heavy-tailed mutations available, it is less critical if an improvement in distance m is missed in phase m . At the same time, since our heavy-tailed mutations also allow to flip more than r bits in phase r , longer phases obtained by taking a larger value of R usually do not have a negative effect on the runtime. For these reasons, the times computed in Theorem 5 depend very little on R . Since the phase length depends only logarithmically on R , we feel that it is safe to choose R as some mildly large constant, say $R = 25$.

The most interesting choice is the value for γ , which sets the balance between the SD-RLS mode of the algorithm and the heavy-tailed mutations. A large rate $1 - \gamma$ of SD-RLS iterations is good to find a single improvement, but can lead to drastic performance losses when there are more improving solutions. Such trade-offs are often to be made in evolutionary computation. For example, the simple RLS heuristic using only 1-bit flips is very efficient on unimodal problems (e.g., has a runtime of $(1 + o(1))n \ln n$ on ONEMAX), but fails on multimodal problems. In contrast, the $(1 + 1)$ EA flips a single bit only with probability approximately $\frac{1}{e}$, and thus optimizes ONEMAX only in time $(1 + o(1))en \ln n$, but can deal with local optima. In a similar vein, a larger value for γ in our algorithm gives some robustness to situations where in phase r other mutations than r -bit flips are profitable – at the price of a slowdown on problems like classic jump functions, where a single improving solution has to be found. It has to be left to the algorithm user to set this trade-off suitably. Taking the example of RLS and the $(1 + 1)$ EA as example, we would generally recommend a constant factor performance loss to buy robustness, that is, a constant value of γ like, e.g., $\gamma = 0.25$.

8 Conclusion

In this work, we proposed a way to combine stagnation detection with heavy-tailed mutation. Our theoretical and experimental results indicate that our new algorithm inherits the good properties of the previous stagnation detection approaches, but is superior in the following respects.

- The additional use of heavy-tailed mutation greatly speed up leaving a local optimum if there is more than one improving solution in a certain distance m . This is because to leave the local optimum, it is not necessary anymore to complete phase $m - 1$.

- Compared to the robust SD-RLS, which is the fairest point of comparison, our algorithm is significantly simpler, as it avoids the two nested loops (implemented via the parameters r and s in [RW21b]) that organize the reversion to smaller rates. Compared to the SD-(1 + 1) EA, our approach can obtain the better runtimes of the SD-RLS approaches in the case that few improving solutions are available, and compared to the simple SD-RLS of [RW21b], our approach surely converges.
- Again comparing our approach to the robust SD-RLS, our approach gives runtimes with exponential tails. Let m be constant. If the robust SD-RLS misses an improvement in distance m in the m -th phase and thus in time $O(n^m)$ – which happens with probability $n^{-\Theta(1)}$ for typical parameter settings –, then strength m is used again only after the $(m + 1)$ -st phase, that is, after $\Omega(n^{m+1})$ iterations. If our algorithm misses such an improvement in phase m , then in each of the subsequent $\ell_{m+1} = \Omega(n^{m+1})$ iterations, it still has a chance of $\frac{\gamma}{2} \binom{n}{m+1}^{-1}$ to find this particular improvement. Hence the probability that finding this improvement takes $\Omega(n^{m+1})$ time, is only $(1 - \frac{\gamma}{2} \binom{n}{m+1}^{-1})^{\Omega(n^{m+1})} \leq \exp(-\frac{\gamma}{2} \Omega(n))$.

As discussed in Section 7, the three parameters of our approach are not too critical to set. For these reasons, we believe that our combination of stagnation detection and heavy-tailed mutation is a very promising approach.

As the previous works on stagnation detection, we have only analyzed stagnation detection in the context of a simple hillclimber. This has the advantage that it is clear that the effects revealed in our analysis are truly caused by our stagnation detection approach. Given that there is now quite some work studying stagnation detection in isolation, for future work it would be interesting to see how well stagnation detection (ideally in the combination with heavy-tailed mutation as proposed in this work) can be integrated into more complex evolutionary algorithms.

Acknowledgement

Amirhossein Rajabi was supported by a research grant by the Danish Council for Independent Research (DFF-FNU 8021-00260B) and a travel grant from the Otto Mønsted foundation.

References

- [ABD20a] Denis Antipov, Maxim Buzdalov, and Benjamin Doerr. Fast mutation in crossover-based algorithms. In *Genetic and Evolutionary Computation Conference, GECCO 2020*, pages 1268–1276. ACM, 2020.
- [ABD20b] Denis Antipov, Maxim Buzdalov, and Benjamin Doerr. First steps towards a runtime analysis when starting with a good solution. In *Parallel Problem Solving From Nature, PPSN 2020, Part II*, pages 560–573. Springer, 2020.
- [ABD21] Denis Antipov, Maxim Buzdalov, and Benjamin Doerr. Lazy parameter tuning and control: choosing all parameters randomly from a power-law distribution. In *Genetic and Evolutionary Computation Conference, GECCO 2021*, pages 1115–1123. ACM, 2021.
- [AD20] Denis Antipov and Benjamin Doerr. Runtime analysis of a heavy-tailed $(1 + (\lambda, \lambda))$ genetic algorithm on jump functions. In *Parallel Problem Solving From Nature, PPSN 2020, Part II*, pages 545–559. Springer, 2020.
- [BBD21] Henry Bambury, Antoine Bultel, and Benjamin Doerr. Generalized jump functions. In *Genetic and Evolutionary Computation Conference, GECCO 2021*, pages 1124–1132. ACM, 2021.
- [COY21] Dogan Corus, Pietro S. Oliveto, and Donya Yazdani. Automatic adaptation of hypermutation rates for multimodal optimisation. In *Foundations of Genetic Algorithms, FOGA 2021*, pages 4:1–4:12. ACM, 2021.
- [DD20] Benjamin Doerr and Carola Doerr. Theory of parameter control for discrete black-box optimization: provable performance gains through dynamic parameter choices. In Benjamin Doerr and Frank Neumann, editors, *Theory of Evolutionary Computation: Recent Developments in Discrete Optimization*, pages 271–321. Springer, 2020. Also available at <https://arxiv.org/abs/1804.05650>.
- [DFK⁺18] Duc-Cuong Dang, Tobias Friedrich, Timo Kötzing, Martin S. Krejca, Per Kristian Lehre, Pietro S. Oliveto, Dirk Sudholt, and Andrew M. Sutton. Escaping local optima using crossover

- with emergent diversity. *IEEE Transactions on Evolutionary Computation*, 22:484–497, 2018.
- [DJW02] Stefan Droste, Thomas Jansen, and Ingo Wegener. On the analysis of the (1+1) evolutionary algorithm. *Theoretical Computer Science*, 276:51–81, 2002.
- [DLMN17] Benjamin Doerr, Huu Phuoc Le, Régis Makhmara, and Ta Duy Nguyen. Fast genetic algorithms. In *Genetic and Evolutionary Computation Conference, GECCO 2017*, pages 777–784. ACM, 2017.
- [Doe20] Benjamin Doerr. Does comma selection help to cope with local optima? In *Genetic and Evolutionary Computation Conference, GECCO 2020*, pages 1304–1313. ACM, 2020.
- [DZ21] Benjamin Doerr and Weijie Zheng. Theoretical analyses of multi-objective evolutionary algorithms on multi-modal objectives. In *Conference on Artificial Intelligence, AAAI 2021*, pages 12293–12301. AAAI Press, 2021.
- [FGQW18a] Tobias Friedrich, Andreas Göbel, Francesco Quinzan, and Markus Wagner. Evolutionary algorithms and submodular functions: Benefits of heavy-tailed mutations. *CoRR*, abs/1805.10902, 2018.
- [FGQW18b] Tobias Friedrich, Andreas Göbel, Francesco Quinzan, and Markus Wagner. Heavy-tailed mutation operators in single-objective combinatorial optimization. In *Parallel Problem Solving from Nature, PPSN 2018, Part I*, pages 134–145. Springer, 2018.
- [FQW18] Tobias Friedrich, Francesco Quinzan, and Markus Wagner. Escaping large deceptive basins of attraction with heavy-tailed mutation operators. In *Genetic and Evolutionary Computation Conference, GECCO 2018*, pages 293–300. ACM, 2018.
- [Prü04] Adam Prügel-Bennett. When a genetic algorithm outperforms hill-climbing. *Theoretical Computer Science*, 320:135–153, 2004.
- [RW20] Amirhossein Rajabi and Carsten Witt. Self-adjusting evolutionary algorithms for multimodal optimization. In *Genetic and*

- Evolutionary Computation Conference, GECCO 2020*, pages 1314–1322. ACM, 2020.
- [RW21a] Amirhossein Rajabi and Carsten Witt. Stagnation detection in highly multimodal fitness landscapes. In *Genetic and Evolutionary Computation Conference, GECCO 2021*, pages 1178–1186. ACM, 2021.
- [RW21b] Amirhossein Rajabi and Carsten Witt. Stagnation detection with randomized local search. In *Evolutionary Computation in Combinatorial Optimization, EvoCOP 2021*, pages 152–168. Springer, 2021.
- [Weg01] Ingo Wegener. Theoretical aspects of evolutionary algorithms. In *Automata, Languages and Programming, ICALP 2001*, pages 64–78. Springer, 2001.
- [Wit21] Carsten Witt. On crossing fitness valleys with majority-vote crossover and estimation-of-distribution algorithms. In *Foundations of Genetic Algorithms, FOGA 2021*, pages 2:1–2:15. ACM, 2021.
- [WQT18] Mengxi Wu, Chao Qian, and Ke Tang. Dynamic mutation based Pareto optimization for subset selection. In *Intelligent Computing Methodologies, ICIC 2018, Part III*, pages 25–35. Springer, 2018.