

Lower bounds for hit-and-run direct search

Jens Jägersküpfer*

Universität Dortmund, Informatik 2, 44221 Dortmund, Germany
JJ@Ls2.cs.uni-dortmund.de

Abstract. “Hit-and-run is fast and fun” to generate a random point in a high dimensional convex set K (Lovász/Vempala, MSR-TR-2003-05). More precisely, the hit-and-run random walk mixes fast independently of where it is started inside the convex set. To hit-and-run from a point $x \in \mathbb{R}^n$, a line L through x is randomly chosen (uniformly over all directions). Subsequently, the walk’s next point is sampled from $L \cap K$ using a membership oracle which tells us whether a point is in K or not.

Here the focus is on black-box optimization, however, where the function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ to be minimized is given as an oracle, namely a black box for f -evaluations. We obtain in an obvious way a direct-search method when we substitute the f -oracle for the K -membership oracle to do a line search over L , and, naturally, we are interested in how fast such a hit-and-run direct-search heuristic converges to the optimum point x^* in the search space \mathbb{R}^n .

We prove that, even under the assumption of perfect line search, the search converges (at best) linearly at an expected rate larger (i. e. worse) than $1 - 1/n$. This implies a lower bound of $0.5n$ on the expected number of line searches necessary to halve the approximation error. Moreover, we show that $0.4n$ line searches suffice to halve the approximation error only with an exponentially small probability of $\exp(-\Omega(n^{1/3}))$. Since each line search requires at least one query to the f -oracle, the lower bounds obtained hold also for the number of f -evaluations.

1 Introduction

Finding an optimum of a given function $f: S \rightarrow \mathbb{R}$ is one of the fundamental problems—in theory as well as in practice. The search space S can be discrete or continuous, like \mathbb{N} or \mathbb{R} . If S has more than one dimension, it may also be a mixture. Here the optimization in “high-dimensional” Euclidean space is considered, i. e., the search space is \mathbb{R}^n . What “high-dimensional” means is usually anything but well defined. A particular 10-dimensional problem in practice may already be considered “high-dimensional” by the one who tries to solve it. Here the crucial aspect is how the optimization time scales with the dimensionality of the search space \mathbb{R}^n , i. e., we consider the optimization time as a function of n . In other words, here we are interested in what happens when the dimensionality

* supported by the German Research Foundation (DFG) through the collaborative research center “Computational Intelligence” (SFB 531)

Copyright 2007 Springer-Verlag Berlin

In: J. Hromkovic et al., SAGA 2007, LNCS 4665, pp. 118-129

of the search space gets higher and higher. This viewpoint is typical for analyses in computer science. In the domain of operations research and mathematical programming, however, focusing on how the optimization time scales with the search space's dimension seems not that common. Usually, the performance of an optimization method is described by means of convergence theory. As an example, let us take a closer look at “Q-linear convergence” (we drop the “Q” in the following): Let \mathbf{x}^* denote the optimum search point of a unimodal function and $\mathbf{x}^{[k]}$ the approximate solution after k optimization steps. Then we have

$$\frac{\text{dist}(\mathbf{x}^*, \mathbf{x}^{[k+1]})}{\text{dist}(\mathbf{x}^*, \mathbf{x}^{[k]})} \rightarrow r \in \mathbb{R}_{<1} \quad \text{as } k \rightarrow \infty$$

where $\text{dist}(\cdot, \cdot)$ denotes some distance measure, most commonly the Euclidean distance between two points (when considering convergence towards \mathbf{x}^* in the search space \mathbb{R}^n , as we do here), or the absolute difference in function value (when considering convergence towards the optimum function value in the objective space). Apparently, there seems to be no connection to n , the dimension of the search space. Yet only if r is an absolute constant, there is actual independence of n . In general, however, the *convergence rate* r depends on n . When we are interested in, say, the number of steps necessary to halve the approximation error (given by the distance from \mathbf{x}^*), the order of this number with respect to n precisely depends on how r depends on n . For instance, if $r = 1 - 0.5/n$, we need $\Theta(n)$ steps; if $r = 1 - 0.5/n^2$, we need $\Theta(n^2)$ steps, and if $r = 1 - 2^{-n}$, we need $2^{\Theta(n)}$ steps. For any fixed dimension, however, in any of the three cases we actually have linear convergence. Thus, the order of convergence tells us something about the “speed” of the optimization, but in general nothing about the n -dependence of the number of steps necessary to ensure a certain approximation error (unless r is an absolute constant, i. e. independent of n). So, in case of linear convergence, we want to know how the convergence rate depends on the dimensionality of the search space.

Methods for solving optimization problems in continuous domains, essentially $S = \mathbb{R}^n$, are usually classified into first-order, second-order, and zeroth-order methods, depending on whether they utilize the gradient (first derivative) of the objective function, the gradient and the Hessian (second derivative), or neither of both. A zeroth-order method is also called *derivative-free* or *direct search*. Newton’s method is a classical second-order method; first-order methods can be (sub)classified into Quasi-Newton, conjugate gradient, and steepest descent methods. Classical zeroth-order methods try to approximate the gradient and to then plug this estimate into a first-order method. Finally, amongst the modern zeroth-order methods, randomized direct-search heuristics like simulated annealing and evolutionary algorithms come into play, which are supposed general-purpose search heuristics.

When information about the gradient is not available, for instance if f relates to a property of some workpiece and is given by computer simulations or even by real-world experiments, then zeroth-order methods are the only option (unless simulations allow for algorithmic/automatic differentiation). As the approximation of the gradient usually involves at least n f -evaluations (forward

finite differences; $2n$ for symmetric finite differences), a single optimization step of a classical zeroth order-method is computationally expensive, in particular if f is given implicitly by complex simulations. In practical optimization, especially in mechanical engineering, this is often the case, and particularly in this field randomized search heuristics (especially evolutionary algorithms) are becoming more and more popular. However, the enthusiasm in practical optimization heuristics has led to an unclear variety of very sophisticated and problem-specific algorithms. Unfortunately, from a theoretical point of view, the development of such algorithms is solely driven by practical success, whereas the aspect of a theoretical analysis is left aside.

In such situations f is given to the optimization algorithm as a black box for f -evaluations (zeroth-order oracle) and the cost of the optimization (the runtime) is defined as the number of queries to this oracle, and we are in the so-called *black-box optimization* scenario. Nemirovsky and Yudin (1983, p. 333) state (w.r.t. optimization in continuous search spaces) in their book *Problem Complexity and Method Efficiency in Optimization*: “From a practical point of view this situation would seem to be more typical. At the same time it is objectively more complicated and it has been studied in a far less extent than the one [with first-order oracles/methods] considered earlier.” After more than two decades there still seems to be some truth in their statement, though to a smaller extent. For discrete black-box optimization, a complexity theory has been successfully started, cf. Droste, Jansen, and Wegener (2006). Lower bounds on the number of f -evaluations (the *black-box complexity*) are proved with respect to classes of functions when an arbitrary(!) optimization heuristic knows about the class \mathcal{F} of functions to which f belongs, but nothing about f itself. The benefits of such results are obvious: They can prove that an allegedly poor performance of an apparently simple black-box algorithm on f is due to \mathcal{F} 's inherent black-box complexity rather than due to the algorithm's simpleness.

As already discussed above, the situation for heuristic optimization in continuous search spaces is different, especially with respect to randomized (direct) methods. The results to be presented here contribute to this emerging field of optimization theory.

2 The Framework for the Randomized Methods

As already noted above, classical zeroth-order methods (i. e. black-box optimizers) for continuous search spaces usually try to approximate the gradient of the function f to be minimized at the current search point \mathbf{x} . Subsequently, a line search along gradient direction is performed to find the next search point, which replaces \mathbf{x} . Usually, the line search aims at locating the best (with respect to the f -value) point on the line through \mathbf{x} , and various strategies for how to do the line search exist (Armijo/Goldstein, Powell/Wolfe, etc.; cf. Nocedal and Wright (2006, Ch. 3) for instance). As the approximation of the gradient usually involves at least n f -evaluations, and as the (approximate) gradient's direction may significantly differ from the direction pointing directly to the optimum \mathbf{x}^*

anyway (cf. ill-conditioned quadratics), more and more direct-search heuristics have been proposed which abandon gradient approximation. Among the first and most prominent ones are the pattern search by Hooke and Jeeves (1961) and the (downhill) simplex method by Nelder and Mead (1965); cf. Kolda, Lewis, and Torczon (2004) for a comprehensive review. Surprisingly, also already in the 1960s randomized direct-search methods were proposed, one is the so-called *evolution strategy* by Rechenberg (1965) and Schwefel (1965). For some obscure reason, however, there has been resentment against randomized algorithms in these early years. This started to change with the randomization of quicksort and randomized testing for primality. At the latest by the time when Dyer, Frieze, and Kannan (1989) came up with a randomized approximation algorithm for the computation of the volume of a convex body in high dimensional space, the (potential) benefits of randomization have won recognition. Though the polynomial expected runtime of this algorithm was not very practical, it showed in principle the power of randomization since for any deterministic algorithm there is a convex set for which the relative approximation error is $n^{\Omega(n)}$ after any polynomial number of steps. At the core of this algorithm was a random walk on a (sufficiently fine) lattice. This algorithm was further improved, in particular by substituting the so-called *ball walk* for the original lattice walk. One step of this ball walk consists in uniformly choosing a point from the hyper-ball of radius δ around the current point. If this point lies in the convex set, then it becomes the next point of the walk. Obviously, one has to choose the parameter δ appropriately. Moreover, when the ball walk is started very close (w. r. t. δ) to the corner of a hypercube, just for instance, it may need an exponential number of steps to leave this corner, making a so-called warm start necessary (i. e. a preprocessing). As recently shown by Lovász and Vempala (2006), using the *hit-and-run walk* instead of the ball walk avoids these two issues. Hit-and-run mixes fast even when started close to the boundary of the convex set, and moreover, no “step size” needs to be appropriately predefined. Also an optimization algorithm based on random walks in convex sets has been proposed (Bertsimas and Vempala, 2004).

As already noted in the abstract, to hit-and-run from a point $\mathbf{x} \in \mathbb{R}^n$ within a convex set $K \subset \mathbb{R}^n$, a line L through \mathbf{x} is randomly chosen (uniformly over all directions). Subsequently, the walk’s next point is sampled from $L \cap K$ (as uniformly as possible) using a membership oracle which tells us whether a sample from L lies in K or not. As also already noted in the abstract, we obtain in an obvious way a hit-and-run direct-search method for black-box optimization of $f: \mathbb{R}^n \rightarrow \mathbb{R}$ when we substitute the f -oracle for the K -membership oracle. Thus, the framework of the heuristics for black-box optimization we consider is as follows: For a given initialization of $\mathbf{x} \in \mathbb{R}^n$ the following loop is performed:

1. Randomly choose a line L through \mathbf{x} (uniformly over all directions).
2. By some kind of a line search (using the f -oracle), find a point $\mathbf{x}' \in L$.
3. Set $\mathbf{x} := \mathbf{x}'$ and GOTO 1 (unless stopping is requested; then output \mathbf{x}).

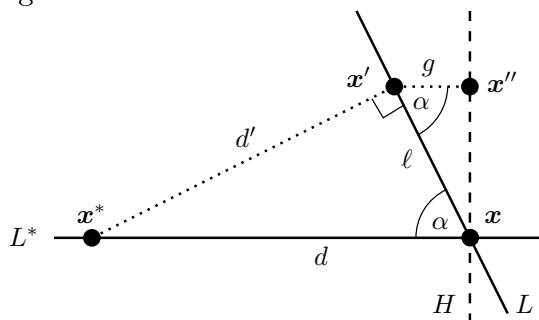
Naturally, we are interested in how fast such a heuristic converges to the optimum point $\mathbf{x}^* \in \mathbb{R}^n$ (we assume that there is a unique global optimum), in particular:

How fast can it converge in principle? That is, we are interested in a general lower bound which is universal for the class of hit-and-run direct-search heuristics.

Note that there are *no* assumptions on how the line search is performed. In particular, for the line search in the i th iteration, the algorithm may use *all* the information gathered from *all* the samples drawn during the preceding $i - 1$ line searches. Naturally, in each step the choice of how to do the line search may additionally depend on the actual direction of L . All in all, a large variety of adaptive strategies for black-box optimization is covered by our framework.

3 The Lower Bounds

Since any reasonable strategy for the line search implies at least one query to the f -oracle, in our scenario the number of f -evaluations is bounded below by the number of line searches. Thus, we focus on the number of line searches in the following and aim at a general lower bound. Therefore, we need an upper bound on the gain of a single line search. We consider the best case: When we want the heuristic to approach the unique optimum point \mathbf{x}^* as fast as possible, we may optimistically assume that \mathbf{x}' was chosen from the line L such that the distance between \mathbf{x}' and \mathbf{x}^* is minimum. Call this a *perfect line search*. The situation is depicted in the figure below.



It is well known that the distance between \mathbf{x}^* and \mathbf{x}' is minimum when $\mathbf{x}' \in L \cap \{\mathbf{x}\}$ is such that the line passing through \mathbf{x}' and \mathbf{x}^* is perpendicular to the line L (given that $\mathbf{x}^* \notin L$, which is the case with probability one, unlike already \mathbf{x} coincides with the optimum point \mathbf{x}^* , because L 's direction is chosen uniformly over all directions).

Let $d := \text{dist}(\mathbf{x}, \mathbf{x}^*)$ denote the current approximation error in the search space and let $d' := \text{dist}(\mathbf{x}', \mathbf{x}^*)$. Furthermore, let L^* denote the line through \mathbf{x} and \mathbf{x}^* . Now consider the hyper-plane H which contains \mathbf{x} and is perpendicular to L^* . Let $\mathbf{x}'' := \arg \min_{\mathbf{y} \in H} \text{dist}(\mathbf{x}', \mathbf{y})$ denote the unique point in H with smallest distance from \mathbf{x}' . Then the angle α between L and L^* equals the angle between L and the line through \mathbf{x}' and \mathbf{x}'' (which is parallel to L^* since it is perpendicular to H just as L^*). Consequently, we have

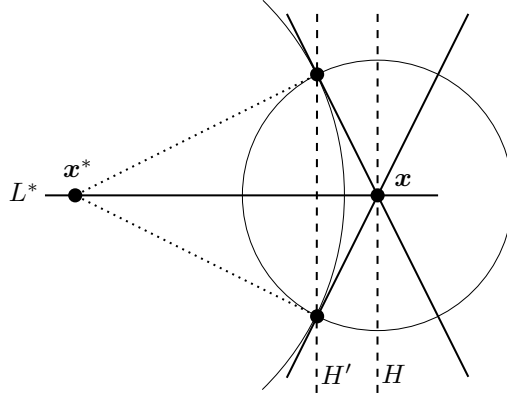
$$d' = d \cdot \sin \alpha \quad \text{and} \quad \text{dist}(\mathbf{x}', H) = \text{dist}(\mathbf{x}', \mathbf{x}) \cdot \cos \alpha.$$

Let $g := \text{dist}(\mathbf{x}', \mathbf{x}'')$ denote the distance of \mathbf{x}' from H , and $\ell := \text{dist}(\mathbf{x}', \mathbf{x})$ so that we have $g/\ell = \cos \alpha$. Since $d'/d = \sin \alpha = \sqrt{1 - (\cos \alpha)^2}$, we obtain

$$\frac{d'}{d} = \sqrt{1 - (g/\ell)^2}, \quad (1)$$

which ranges in $[0, 1]$ since $g \in [0, \ell]$. Thus, instead of focusing on the distribution of $\sin \alpha$ when the line L is chosen uniformly over all directions, we can focus on the ratio g/ℓ and concentrate on the distribution of this *relative distance* of \mathbf{x}' from the hyper-plane H (namely, relative to the distance of \mathbf{x}' from \mathbf{x}). (It will shortly become clear why this makes sense.)

In two dimensions, like in the figure above, for any fixed $d' \in (0, d)$ there are exactly two (different) lines through \mathbf{x} with distance d' from the optimum point \mathbf{x}^* . (Note that by fixing d' we also fixed ℓ and g .) In three or more dimensions, however, there is an infinite number of such lines. In three dimensions they form a double cone with its apex at \mathbf{x} , and all points of this cone with a distance of exactly d' from \mathbf{x}^* (namely all \mathbf{x}') form a circle. This circle lies in a plane which is parallel to H (a plane in three dimensions). In general, i.e. in $n \geq 3$ dimensions, the potential points \mathbf{x}' form the set $S := \{\mathbf{x}' \in \mathbb{R}^n \mid \text{dist}(\mathbf{x}', \mathbf{x}^*) = d' \text{ and } \text{dist}(\mathbf{x}', \mathbf{x}) = \ell\}$, which is an $(n-1)$ -sphere since S is the intersection of two hyper-spheres, namely of the hyper-sphere with radius d' centered at \mathbf{x}^* and the hyper-sphere with radius ℓ centered at \mathbf{x} . Moreover, S lies in the hyper-plane H' which is parallel to H such that it has distance g from H and distance $d - g$ from \mathbf{x}^* . The situation is depicted below, where the left sphere consists of all points with distance d' from the optimum point \mathbf{x}^* , and the right sphere consists of all points with distance ℓ from our current approximate solution \mathbf{x} .



Recall that we fixed $d' \in (0, d)$ for the above discussion, and that this implies fixed values for ℓ and $g = \text{dist}(H', H)$. Now consider a randomly chosen line L through \mathbf{x} (uniform over all directions). According to our construction, if L penetrates the $(n-1)$ -sphere $S \subset H'$, then the perfect line search on L yields a point with a distance of exactly d' from \mathbf{x}^* . Now, if L lies inside the double cone, i.e., L penetrates the open $(n-1)$ -ball whose missing boundary is S , then the perfect line search yields a point with a distance smaller than d' from \mathbf{x}^* . If L

lies outside the double cone (except for passing through the apex \mathbf{x} , of course), then the perfect line search yields a point with a distance larger than d' from \mathbf{x}^* . Thus, we are interested in the probability that L is chosen such that it lies inside the cone, namely the probability that the perfect line search yields a point with a distance of less than d' from \mathbf{x}^* .

Now, how can we actually pick a line through \mathbf{x} such that its direction is uniformly random? We pick uniformly at random a point \mathbf{y} from/over the unit hyper-sphere centered at \mathbf{x} and choose L as the line through \mathbf{y} and \mathbf{x} . From this point of view, the perfect line search yields a point with a distance of exactly d' from \mathbf{x}^* if \mathbf{y} 's distance from H is exactly g/ℓ ; a point with a distance smaller than d' from \mathbf{x}^* if \mathbf{y} 's distance from H is larger than g/ℓ ; and a point with a distance larger than d' from \mathbf{x}^* if \mathbf{y} 's distance from H is smaller than g/ℓ .

In other words, we can consider the random variable $R := d'/d$ as a function of the random variable G defined as \mathbf{y} 's distance from the hyper-plane H , where the point \mathbf{y} is chosen uniformly over the unit hyper-sphere centered at \mathbf{x} . Namely, we have $R = \sqrt{1 - G^2}$, cf. Equation 1. (Note that the distribution of \mathbf{y} over \mathbb{R}^n is spherically symmetric; more precisely, invariant w.r.t. orthonormal transformations.) For $n \geq 4$ the density function of G 's distribution over $[0, 1]$ is given by $(1 - x^2)^{(n-3)/2} / \Psi$ (Jägersküpfer, 2003), where $\Psi = \int_0^1 (1 - x^2)^{(n-3)/2} dx$ (normalization) and the value of this integral is $\Psi = \sqrt{\pi/4} \cdot \Gamma(n/2 - 1/2) / \Gamma(n/2) = \sqrt{\pi/n}/2 + \Theta(n^{-3/2})$, where “ Γ ” denotes the well-known gamma function. Consequently, \mathbf{y} 's expected distance from H equals $\int_0^1 x \cdot (1 - x^2)^{(n-3)/2} dx / \Psi = (n - 1)^{-1} / \Psi$ which turns out to be about $0.8/\sqrt{n}$. This might appear bewildering (at first) since this implies that, as the search space's dimensionality increases, the expected distance from H tends to zero—although \mathbf{y} 's distance from \mathbf{x} is fixed to one and H is hit with zero probability. However, noting that H is an affine subspace with dimension $n - 1$ (i.e. codimension 1), it may become more plausible that getting far away from H becomes less and less probable as n increases. It might help even more to recall that an n -hypercube with unit diameter (longest diagonal) has edges of length $1/\sqrt{n}$.

Naturally, $\mathbb{E}[G]$ does not tell us much about $\mathbb{E}[R] = \mathbb{E}[\sqrt{1 - G^2}]$, the expectation in which we are actually interested. We can easily compute it, though:

$$\mathbb{E}[R] = \int_0^1 \sqrt{1 - x^2} \cdot (1 - x^2)^{(n-3)/2} / \Psi dx = \int_0^1 (1 - x^2)^{n/2-1} dx / \Psi.$$

Since $\int_0^1 (1 - x^2)^{n/2-1} dx = \sqrt{\pi/4} \cdot \Gamma(n/2) / \Gamma(n/2 + 1/2)$, we obtain

$$\mathbb{E}[R] = \frac{\sqrt{\pi/4} \cdot \Gamma(n/2) / \Gamma(n/2 + 1/2)}{\sqrt{\pi/4} \cdot \Gamma(n/2 - 1/2) / \Gamma(n/2)} = \frac{\Gamma(n/2) \cdot \Gamma(n/2)}{\Gamma(n/2 - 1/2) \cdot \Gamma(n/2 + 1/2)}. \quad (2)$$

Using $\Gamma(n/2 + 1/2) = \Gamma(n/2 - 1/2) \cdot (n/2 - 1/2)$, we have

$$\mathbb{E}[R] = \frac{n - 1}{2} \cdot \left(\frac{\Gamma(n/2)}{\Gamma(n/2 + 1/2)} \right)^2,$$

and since $\Gamma(n/2 + 1/2)/\Gamma(n/2) < \sqrt{n/2}$, we obtain the following lower bound on the expected factor by which the approximation error is reduced in each step:

$$\mathbb{E}[R] > \frac{n-1}{2} \cdot \frac{2}{n} = 1 - \frac{1}{n}.$$

This lower bound holds for perfect line search and, as a direct consequence, also for any other line-search strategy. With other words, this bound is universal for the class of hit-and-run direct-search methods.

To see how good this general lower bound on $\mathbb{E}[R]$ actually is, an upper bound on $\mathbb{E}[R]$ under the assumption of perfect line search would be nice. Using that $\Gamma(n) = (n-1)!$, $\Gamma(n/2) = (n-2)!! \cdot \sqrt{\pi}/2^{(n-1)/2}$, and $\Gamma(k+1/2) = (2k-1)!! \cdot \sqrt{\pi}/2^k$ (where $k!!$ is defined as $2 \cdot 4 \cdot 6 \cdots k$ for even k , and as $1 \cdot 3 \cdot 5 \cdots k$ for odd k), the right-hand side of Equation 2 can be estimated as follows:

$$\mathbb{E}[R] < \frac{2n-1}{2n} = 1 - \frac{1}{2n}.$$

In other words, for perfect line search, the expected factor by which the approximation error is reduced (in each step) is smaller (i.e. better) than $1 - 0.5/n$. This shows that our general lower bound of $1 - 1/n$ on $\mathbb{E}[R]$ is actually pretty tight. All in all, we have proved the following result:

Theorem 1. *Consider the optimization of a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ with a unique optimum point \mathbf{x}^* . Then we have for $n \geq 4$:*

The (hypothetical) hit-and-run direct-search method which performs a perfect line search in each step converges linearly towards \mathbf{x}^ at an expected rate of $1 - \beta/n$, where $0.5 < \beta < 1$ (and β depends on n according to Equation 2).*

Independently of how a hit-and-run direct-search method performs the line searches, the expected factor by which the approximation error (i.e. the distance from \mathbf{x}^) is reduced is larger than $1 - 1/n$ in each step. That is, if (at all) a hit-and-run direct-search method converges towards \mathbf{x}^* , then at best linearly at an expected rate larger (i.e. worse) than $1 - 1/n$.*

The result on the expected factor by which the approximation error is reduced directly implies a bound on the expected spatial gain towards the optimum point \mathbf{x}^* . Therefore, let $d^{[i]}$ denote the approximation error (i.e. the distance from \mathbf{x}^*) after the i th step, and let $d^{[0]}$ denote the initial approximation error. For a fixed $d^{[i-1]}$, let $\Delta^{[i]} := d^{[i]} - d^{[i-1]}$ be defined as the random variable corresponding to the spatial gain towards \mathbf{x}^* in the i th step. Then the above theorem says that in general, i.e. for any hit-and-run direct-search method, $\mathbb{E}[\Delta^{[i]}] < d^{[i-1]}/n$ in each step i . Moreover, for perfect line search, in each step $\mathbb{E}[\Delta^{[i]}] = \beta(n) \cdot d^{[i-1]}/n$ for some function $\beta: \mathbb{N} \rightarrow (0.5, 1)$.

Let us stick with perfect line search in the following. Then the approximation error is non-increasing, i.e., $d^{[0]} \geq d^{[1]} \geq d^{[2]} \dots$ (actually, $d^{[i+1]} < d^{[i]}$ with probability one, since the randomly chosen line lies in H with zero probability). Thus, in each step $\Delta^{[i]} < d^{[i-1]}/n \leq d^{[0]}/n$, and consequently, the number of steps necessary for an expected total gain of at least $d^{[0]}/2$ is larger than

$(d^{[0]}/2)/(d^{[0]}/n) = n/2$ (actually, $n \ln 2$). However, in general, maximizing the *expected total gain* of a fixed number of steps need not necessarily result in minimizing the *expected number of steps* to realize a specified gain. Nevertheless, $n/2$ will turn out to be a lower bound on the expected number of steps which are necessary to halve the approximation error. The proof is a straight-forward application of the following lemma, which is a modification of Wald’s equation.

Lemma 2. *Let X_1, X_2, \dots denote random variables with bounded range and S the random variable defined by $S = \min\{t \mid X_1 + \dots + X_t \geq g\}$ for a given $g > 0$. Given that S is a stopping time (i. e., the event $\{S = t\}$ depends only on X_1, \dots, X_t), if $E[S] < \infty$ and $E[X_i \mid S \geq i] \leq u \neq 0$ for all X_i , then $E[S] \geq g/u$.*

(A proof can be found, e. g., in Jägersküpfer, 2007.) We let X_i denote $\Delta^{[i]}$ and choose $g := d^{[0]}/2$. As we have just seen, $0 \leq \Delta^{[i]} \leq d^{[0]}$, and since in our scenario “ $S \geq i$ ” merely means that the approximation error has not been halved in the first $i-1$ steps, actually $E[\Delta^{[i]} \mid S \geq i] < d^{[0]}/n =: u$. Finally, we note that S is in fact a stopping time so that $g/u = n/2$ is indeed a lower bound on the expected number of steps to halve the approximation error (unless $E[S]$ was infinite, in which case we would not need to prove a lower bound anyway). Due to the linearity of expectation, the expected number of steps to halve the approximation error $b \in \mathbb{N}$ times is lower bounded by $(n/2) + (b-1) \cdot (n/2 - 1)$, where the rightmost “ -1 ” emerges because the last step within a halving-phase is also (and must be counted as) the first step of the following halving-phase. Thus, we have just proved the following result.

Theorem 3. *Let a hit-and-run direct-search method optimize a function in \mathbb{R}^n , $n \geq 4$, with a unique optimum. Let $b: \mathbb{N} \rightarrow \mathbb{N}$. For perfect line search, the expected number of steps until the approximation error in the search space is less than a $2^{-b(n)}$ -fraction of the initial one is lower bounded by $b(n) \cdot n/2 - b(n) + 1$.*

Now that we know that at least $n/2$ steps are necessary *in expectation* to halve the approximation error, we would like to know whether there is a good chance of getting by with considerably fewer steps. In fact, we want to show that there is almost no chance of getting by with a little fewer steps. Actually, we are going to prove that $0.4n$ steps suffice to halve the approximation error only with an exponentially small probability. Therefore recall the following notions and notations, where X and Y denote random variables:

- X *stochastically dominates* Y , in short “ $X \succ Y$,” if (and only if) $\forall a \in \mathbb{R}: \mathbb{P}\{X \leq a\} \leq \mathbb{P}\{Y \leq a\}$. Obviously, “ \succ ” is a transitive relation.
- If $X \succ Y$ as well as $Y \succ X$, i. e., $\forall a \in \mathbb{R}: \mathbb{P}\{X \leq a\} = \mathbb{P}\{Y \leq a\}$, then X and Y are equidistributed and we write “ $X \sim Y$.”

Theorem 4. *Let a hit-and-run direct-search method optimize a function in \mathbb{R}^n with a unique optimum. Let $b: \mathbb{N} \rightarrow \mathbb{N}$ such that $b(n) = \text{poly}(n)$. For perfect line search, with a very high probability of $1 - \exp(-\Omega(n^{1/3}))$ more than $b(n) \cdot 0.4n$ steps are necessary until the approximation error is less than a $2^{-b(n)}$ -fraction of the initial approximation error.*

Proof. Assume that $\mathbf{x}^{[0]} \neq \mathbf{x}^*$. Because in each step perfect line search is performed, $\Delta^{[i]}/d^{[i-1]} \sim \Delta^{[j]}/d^{[j-1]}$ for $i, j \in \mathbb{N}$ (scale invariance). Since moreover $d^{[0]} \geq d^{[1]} \geq d^{[2]} \dots$, we have $\Delta^{[1]} \succ \Delta^{[2]} \succ \dots$ for the single-step gains. Let X_1, X_2, X_3, \dots denote independent instances of the random variable $\Delta^{[1]}$. Then $\forall i \in \mathbb{N}: X_i \succ \Delta^{[i]}$, and hence $\sum_{i=1}^k \Delta^{[i]} \prec S_k := \sum_{i=1}^k X_i$. In less formal words: Adding up k independent instances of the random variable which corresponds the spatial gain in the first step results in a random variable (namely S_k) which stochastically dominates the random variable given by the total gain of the first k steps. The advantage of considering S_k instead of the “true” total gain of these steps is the following: S_k is the sum of independent random variables so that we can apply Hoeffding’s bound. Namely, Hoeffding (1963, Theorem 2) tells us:

Let X_1, \dots, X_k denote independent random variables with bounded ranges so that $a_i \leq X_i \leq b_i$ with $a_i < b_i$ for $i \in \{1, \dots, k\}$. Let $S := X_1 + \dots + X_k$. Then $\mathbb{P}\{S \geq \mathbb{E}[S] + x\} \leq \exp(-2x^2 / \sum_{i=1}^k (b_i - a_i)^2)$ for any $x > 0$.

If the support of each random variable X_i is contained in $[a, b] \subset \mathbb{R}$, Hoeffding’s bound becomes $\exp(-2 \cdot (x/(b-a))^2/k)$. So, let $k := 0.4n$ and $S := S_k$. Then $\mathbb{E}[S] = 0.4n \cdot \mathbb{E}[\Delta^{[1]}] \leq 0.4d^{[0]}$, and for the application of Hoeffding’s bound we choose $x := 0.1d^{[0]}$, which yields an upper bound of $\exp(-0.05(d^{[0]}/(b-a))^2/n)$ on the probability that the approximation error is halved in $0.4n$ steps. We can choose $a := 0$ so that we obtain $\mathbb{P}\{X_1 + \dots + X_k \geq d^{[0]}/2 \mid X_1, \dots, X_k \leq b\} \leq \exp(-0.05(d^{[0]}/b)^2/n)$, where b is an upper bound on the gain towards the optimum point \mathbf{x}^* in a step. Unfortunately, when substituting the trivial upper bound of $d^{[0]}$ for b , the upper bound on the probability becomes $\exp(-0.05/n)$, which tends to one as n grows. For $b := d^{[0]}/n^{2/3}$, however, we obtain (recall that k was chosen as $0.4n$)

$$\mathbb{P}\left\{X_1 + \dots + X_k \geq d^{[0]}/2 \mid X_1, \dots, X_k \leq d^{[0]}/n^{2/3}\right\} \leq e^{-0.05n^{1/3}}.$$

Thus, if we can show that $\mathbb{P}\{X_i > d^{[0]}/n^{2/3}\} = e^{-\Omega(n^{1/3})}$ in each of the $0.4n$ steps, we obtain (by an application of the union bound)

$$\mathbb{P}\left\{X_1 + \dots + X_k \geq d^{[0]}/2\right\} \leq e^{-0.05n^{1/3}} + 0.4n \cdot e^{-\Omega(n^{1/3})} = e^{-\Omega(n^{1/3})}.$$

Finally, by another application of the union bound, we obtain the theorem because $b(n) = \text{poly}(n)$ implies $b(n) \cdot e^{-\Omega(n^{1/3})} = e^{-\Omega(n^{1/3})}$.

In other words, it remains to be shown that $\mathbb{P}\{\Delta^{[0]} > d^{[0]}/n^{2/3}\}$ is actually bounded above by $e^{-\Omega(n^{1/3})}$. Therefore, recall Equation 1. It tells us that $d - d' = d \cdot (1 - \sqrt{1 - (g/\ell)^2})$. As a consequence, $\mathbb{P}\{\Delta > d/n^{2/3}\}$ is equal to $\mathbb{P}\{1 - \sqrt{1 - G^2} > 1/n^{2/3}\}$. Solving the inequality $1 - \sqrt{1 - G^2} > 1/n^{2/3}$ for G yields $G > \sqrt{2/n^{2/3} + 1/n^{4/3}}$ so that that $\Delta > d/n^{2/3}$ actually implies $G > \sqrt{2}/n^{1/3}$. Since G ’s density is a non-increasing function in $[0, 1]$,

$$\mathbb{P}\left\{G > \sqrt{2}/n^{1/3}\right\} = \int_{\sqrt{2}/n^{1/3}}^1 (1 - x^2)^{(n-3)/2} dx < \left(1 - \frac{2}{n^{2/3}}\right)^{(n-3)/2}.$$

Since $(1-t/k)^k < e^{-t}$ for $0 < t < k > 1$, we have $(1 - 2/n^{2/3})^{n^{2/3}} < e^{-2}$, so that $\mathbb{P}\{\Delta > d/n^{2/3}\} < \mathbb{P}\{G > \sqrt{2}/n^{1/3}\} < e^{-2 \cdot ((n-3)/2)/n^{2/3}} = e^{-n^{1/3} + 3/n^{2/3}}$. \square

4 Discussion and Conclusion

Even though it is clear from intuition that the lower bounds presented in the two preceding theorems do not only hold for perfect line search but for any line-search strategy, they are formally proved only for perfect line search. Interestingly, we can easily show that our theorems do hold independently of how the line searching is actually done: By induction over the number of steps i we show that the random variable which corresponds to the approximation error after i steps for a given line-search strategy stochastically dominates the random variable $d^{[i]}$ for perfect line search, which we considered in the proofs.

So, hit-and-run direct-search methods converge (at best and if at all) linearly with an expected rate larger/worse than $1 - 1/n$. In simple words, the reason for this is that in high dimensions the randomly chosen direction is with a high probability “almost perpendicular” to the direction pointing directly towards the optimum point \mathbf{x}^* . For the further discussion, consider the simple toy problem of minimizing a quadratic form $\mathbf{x} \mapsto \mathbf{x}^\top \mathbf{Q} \mathbf{x}$, where the matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is positive definite. For this simple scenario, steepest descent converges at least linearly at a rate which is independent of the dimension n but which gets worse when the condition number of \mathbf{Q} increases—when assuming a worst-case starting point (cf. Nocedal and Wright (2006, Sec. 3.3) for instance). In the best case, however, steepest descent needs a single (perfect) line search to determine the optimum. Thus, for ill-conditioned quadratics, the performance of steepest descent heavily depends on the starting point. This is one reason why usually preconditioning is applied. Hypothetically assume for a moment the extreme of perfect preconditioning, so that $\mathbf{x}^\top \mathbf{I} \mathbf{x} = |\mathbf{x}|^2$ is to be minimized. Interestingly, the original evolution strategy from 1965 by Rechenberg/Schwefel mentioned in the introduction, a very simple randomized method which belongs to the class of hit-and-run direct-search methods (a line search consists in sampling a single point), actually gets by with $O(n)$ f -evaluations with very high probability to halve the approximation error in this scenario (Jägersküpper, 2003). This shows that the very general lower bound obtained here can be met at least up to a constant factor. However, in this ideal scenario steepest descent needs a single (perfect) line search to find the optimum independently of the starting point. Now, as we consider black-box optimization, steepest descent must approximate the gradient. Even though the approximation of the gradient may cost $2n$ f -evaluations, a single line search in this approximate direction may yield a significantly larger gain towards the optimum—whereas a hit-and-run method needs at least $0.5n$ f -evaluations to halve the approximation error *in any case* (in expectation; $0.4n$ with very high probability). Thus, with a passable preconditioning, the approximation of the gradient should pay off—even though it costs a linear (in n) number of f -evaluations per step—so that it will likely be superior to hit-and-run into a random direction.

As just discussed, for smooth functions we cannot expect hit-and-run direct search to compete with methods which learn (and then utilize) second-order information like the well-known BFGS method or (nonlinear) conjugate gradient methods. Clearly, hit-and-run *can* make sense in real-world optimization—when classical/established (nonlinear) methods have turned out to fail. For instance, when the function to be optimized is non-smooth, or disturbed by noise, or highly multimodal such that gradient approximation is deceptive. Then, however, as we have proved here, we should not expect such hit-and-run direct search to be fast.

References

- Bertsimas, D., Vempala, S. (2004): *Solving convex programs by random walks*. Journal of the ACM, 51(4):540–556.
- Droste, S., Jansen, T., Wegener, I. (2006): *Upper and lower bounds for randomized search heuristics in black-box optimization*. Theory of Computing Systems, 39(4):525–544.
- Dyer, M. E., Frieze, A. M., Kannan, R. (1989): *A random polynomial time algorithm for approximating the volume of convex bodies*. In *Proceedings of the 21st Annual ACM Symposium on Theory of Computing (STOC)*, 375–381.
- Fogel, D. B. (editor) (1998): *Evolutionary Computation: The Fossil Record*. Wiley-IEEE Press.
- Hoeffding, W. (1963): *Probability inequalities for sums of bounded random variables*. American Statistical Association Journal, 58(301):13–30.
- Hooke, R., Jeeves, T. A. (1961): *“Direct search” solution of numerical and statistical problems*. Journal of the ACM, 8(2):212–229.
- Jägersküpfer, J. (2003): *Analysis of a simple evolutionary algorithm for minimization in Euclidean spaces*. In *Proc. 30th Int’l Colloquium on Automata, Languages and Programming (ICALP)*, vol. 2719 of LNCS, 1068–79, Springer.
- Jägersküpfer, J. (2007): *Algorithmic analysis of a basic evolutionary algorithm for continuous optimization*. Theoretical Computer Science, 379(3):329–347.
- Kolda, T. G., Lewis, R. M., Torczon, V. (2004): *Optimization by direct search: New perspectives on some classical and modern methods*. SIAM Review, 45(3):385–482.
- Lovász, L., Vempala, S. (2006): *Hit-and-run from a corner*. SIAM Journal on Computing, 35(4):985–1005.
- Nelder, J. A., Mead, R. (1965): *A simplex method for function minimization*. The Computer Journal, 7:308–313.
- Nemirovsky, A. S., Yudin, D. B. (1983): *Problem Complexity and Method Efficiency in Optimization*. Wiley, New York.
- Nocedal, J., Wright, S. (2006): *Numerical Optimization*. Springer, 2nd edn.
- Rechenberg, I. (1965): *Cybernetic solution path of an experimental problem*. Royal Aircraft Establishment, in Fogel (1998).
- Schwefel, H.-P. (1965): *Kybernetische Evolution als Strategie der experimentellen Forschung in der Strömungstechnik*. Diploma thesis, Technische Universität Berlin.