# Global convergence rate analysis of unconstrained optimization methods based on probabilistic models

C. Cartis[*]      K. Scheinberg[†]

January 6, 2017

## Abstract

We present global convergence rates for a line-search method which is based on random first-order models and directions whose quality is ensured only with certain probability. We show that in terms of the order of the accuracy, the evaluation complexity of such a method is the same as its counterparts that use deterministic accurate models; the use of probabilistic models only increases the complexity by a constant, which depends on the probability of the models being good. We particularize and improve these results in the convex and strongly convex case.

We also analyze a probabilistic cubic regularization variant that allows approximate probabilistic second-order models and show improved complexity bounds compared to probabilistic first-order methods; again, as a function of the accuracy, the probabilistic cubic regularization bounds are of the same (optimal) order as for the deterministic case.

**Keywords:** line-search methods, cubic regularization methods, random models, global convergence analysis.

## 1 Introduction

We consider in this paper the unconstrained optimization problem

$$\min_{x \in \mathbb{R}^n} f(x),$$

where the first (and second, when specified) derivatives of the objective function $f(x)$ are assumed to exist and be (globally) Lipschitz continuous.

Most unconstrained optimization methods rely on approximate local information to compute a local descent step in such a way that sufficient decrease of the objective function is achieved.

To ensure such sufficient decrease, the step has to satisfy certain requirements. Often in practical applications ensuring these requirements for each step is prohibitively expensive or impossible. This may be due to the fact that derivative information about the objective function is not available or because full gradient (and Hessian) are too expensive to compute, or a model of the objective function is too expensive to optimize accurately.

Recently, there has been a significant increase in interest in unconstrained optimization methods with inexact information. Some of these methods consider the case when gradient information is inaccurate. This error in the gradient computation may simply be bounded in the worst case (deterministically), see, for example, [11, 20], or the error is random and the estimated gradient is accurate in expectation, as in stochastic gradient algorithms, see for example, [12, 19, 23, 21]. These methods are typically applied in a convex setting and do not extend to nonconex cases. Complexity bounds are derived that bound the expected accuracy that is achieved after a given number of iterations.

In the nonlinear optimization setting, the complexity of various unconstrained methods has been derived under exact derivative information [7, 8, 17], and also under inexact information, where the errors are bounded in a deterministic fashion [3, 6, 11, 14, 20]. In all the cases of the deterministic inexact setting, traditional optimization algorithms such as line search, trust region or adaptive regularization algorithms are applied with little modification and work in practice as well as in theory, while the error is assumed to be bounded in some decaying manner *at each iteration*. In contrast, the methods based on stochastic estimates of the derivatives, do not assume deterministically bounded errors, however they are quite different from the "traditional" methods in their strategy for step size selection and averaging of the iterates. In other words, they are not simple counterparts of the deterministic methods.

Our purpose in this paper is to derive a class of methods which inherit the best properties of traditional deterministic algorithms, and yet relax the assumption that the derivative/model error is bounded in a deterministic manner. Moreover, we do not assume that the error is zero in expectation or that it has a bounded variance. Our results apply in the setting where at each iteration, with sufficiently high probability, the error is bounded in a decaying manner, while in the remaining cases, this error can be arbitrarily large. In this paper, we assume that the error may happen in the computation of the derivatives and search directions, but that there is no error in the function evaluations, when success of an iterate has to be validated.

Recently several methods for unconstrained black-box optimization have been proposed, which rely on random models or directions [1, 13, 16], but are applied to deterministic functions. In this paper we take this line of work one step further by establishing expected convergence rates for several schemes based on one generic analytical framework.

We consider four cases and derive four different complexity bounds. In particular, we analyze a line search method based on random models, for the cases of general nonconex, convex and strongly convex functions. We also analyze a second order method - an adaptive regularization method with cubics [7, 8] - which is known to achieve the optimal convergence rate for the nonconvex smooth functions [5] and we show that the same convergence rate holds in expectation.

In summary, our results differ from existing literature using inexact, stochastic or random information in the following main points:

- Our models are assumed to be "good" with some probability, but there is no other assumptions on the expected values or variance of the model parameters.

- The methods that we analyze are essentially the exact counterparts of the deterministic

methods, and do not require averaging of the iterates or any other significant changes. We believe that, amongst other things, our analysis helps to understand the convergence properties of practical algorithms, that do not always seek to ensure theoretically required model quality.

- Our main convergence rate results provide a bound on the *expected number of iterations* that the algorithms take before they achieve a desired level of accuracy. This is in contrast to a typical analysis of randomized or stochastic methods, where what is bounded is the expected error after a given number of iterations. Both bounds are useful, but we believe that the bound on the expected number of steps is a somewhat more meaningful complexity bound in our setting. The only other work that we are aware of which provides bounds in terms of the number of required steps is [13] where probabilistic bounds are derived in the particular context of random direct search with possible extension to trust region methods as discussed in Section 6 of [13].

An additional goal of this paper is to present a general theoretical framework, which could be used to analyze the behavior of other algorithms, and different possible model construction mechanisms under the assumption that the objective function is deterministic. We propose a general analysis of an optimization scheme by reducing it to the analysis of a stochastic process. Convergence results for a trust region method in [1] also rely on a stochastic process analysis, but only in terms of behavior in the limit. These results have now been extended to noisy (stochastic) functions, see [9, 10]. Deriving convergence *rates* for methods applied to stochastic functions is the subject of future work and is likely to depend on the results in this paper.

The rest of the paper is organized as follows. In Section 2 we describe the general scheme which encompasses several unconstrained optimization methods. This scheme is based on using random models, which are assumed to satisfy some "quality" conditions with probability at least $p$, conditioned on the past. Applying this optimization scheme results in a stochastic process, whose behavior is analyzed in the later parts of Section 2. Analysis of the stochastic process allows us to bound the expected number of steps of our generic scheme until a desired accuracy is reached. In Section 3 we analyze a linesearch algorithm based on random models and show how its behavior fits into our general framework for the cases of nonconex, convex and strongly convex functions. In Section 4 we apply our generic analysis to the case of the Adaptive Regularization method with Cubics (ARC). Finally, in Section 5 we describe different settings where the models of the objective functions satisfy the probabilistic conditions of our schemes.

## 2   A general optimization scheme with random models

This section presents the main features of our algorithms and analysis, in a general framework that we will, in subsequent sections, particularize to specific algorithms (such as linesearch and cubic regularization) and classes of functions (convex, nonconvex). The reasons for the initial generic approach is to avoid repetition of the common elements of the analysis for the different algorithms and to emphasize the key ingredients of our analysis, which is possibly applicable to other algorithms (provided they satisfy our framework).

## 2.1 A general optimization scheme

We first describe a generic algorithmic framework that encompasses the main components of the unconstrained optimization schemes we analyze in this paper. The scheme relies on building a model of the objective function at each iteration, minimizing this model or reducing it in a sufficient manner and considering the step which is dependent on a stepsize parameter and which provides the model reduction (the stepsize parameter may be present in the model or independent of it). This step determines a new candidate point. The function value is then computed (accurately) at the new candidate point. If the function reduction provided by the candidate point is deemed sufficient, then the iteration is declared successful, the candidate point becomes the new iterate and the step size parameter is increased. Otherwise, the iteration is unsuccessful, the iterate is not updated and the step size parameter is reduced.

We summarize the main steps of the scheme below.

**Algorithm 2.1 Generic optimization framework based on random models**

**Initialization**

> *Choose a class of (possibly random) models $m_k(x)$, choose constants $\gamma \in (0,1)$, $\theta \in (0,1)$, $\alpha_{\max} > 0$. Initialize the algorithm by choosing $x_0$, $m_0(x)$, $0 < \alpha_0 < \alpha_{\max}$.*

**1. Compute a model and a step**

> *Compute a local (possibly random) model $m_k(x)$ of $f$ around $x^k$.*
> *Compute a step $s^k(\alpha_k)$ which reduces $m_k(x)$, where the parameter $\alpha_k > 0$ is present in the model or in the step calculation.*

**2. Check sufficient decrease**

> *Compute $f(x^k + s^k(\alpha_k))$ and check if sufficient reduction (parametrized by $\theta$) is achieved in $f$ with respect to $m_k(x^k) - m_k(x^k + s^k(\alpha_k))$.*

**3. Successful step**

> *If sufficient reduction is achieved then, $x^{k+1} := x^k + s^k(\alpha_k)$, set $\alpha_{k+1} = \min\{\alpha_{\max}, \gamma^{-1}\alpha_k\}$. Let $k := k + 1$.*

**4. Unsuccessful step**

> *Otherwise, $x^{k+1} := x^k$, set $\alpha_{k+1} = \gamma\alpha_k$. Let $k := k + 1$.*

Let us illustrate how the above scheme relates to standard optimization methods. In *linesearch methods*, one minimizes a linear model $m_k(x) = f(x^k) + (x - x^k)^T g^k$ (subject to some normalization), or a quadratic one $m_k(x) = f(x^k) + (x - x^k)^T g^k + \frac{1}{2}(x - x^k)^\top b^k(x - x^k)$ (when the latter is well-defined, with $b^k$ - a Hessian approximation matrix), to find directions $d^k = -g^k$ or $d^k = -(b^k)^{-1}g^k$, respectively. Then the step is defined as $s^k(\alpha_k) = \alpha_k d^k$ for some $\alpha_k$ and, commonly, the (Armijo) decrease condition is checked,

$$f(x^k) - f(x^k + s^k(\alpha_k)) \geq -\theta s^k(\alpha_k)^T g^k,$$

where $-\theta s^k(\alpha_k)^T g^k$ is a multiple of $m_k(x^k) - m_k(x^k + s^k(\alpha_k))$. Note that if the model stays the same in that $m_k(x) \equiv m_{k-1}(x)$ for each $k$, such that $(k-1)$st iteration is unsuccessful, then the above framework essentially reduces to a standard deterministic linesearch.

In the case of *cubic regularization methods*, $s^k(\alpha_k)$ is computed to approximately minimize a cubic model $m_k(x) = f(x^k) + (x - x^k)^T g^k + \frac{1}{2}(x - x^k)^\top b^k (x - x^k) + \frac{1}{3\alpha_k}\|x - x^k\|^3$ and the sufficient decrease condition is

$$\frac{f(x^k) - f(x^k + s^k(\alpha_k))}{m(x^k) - m(x^k + s^k(\alpha_k))} \geq \theta > 0.$$

Note that here as well, in the deterministic case, $g^k = g^{k-1}$ and $b^k = b^{k-1}$ for each $k$ such that $(k-1)$st iteration is unsuccessful but $\alpha_k \neq \alpha_{k-1}$.

The key assumption in the usual deterministic case is that the models $m_k(x)$ are sufficiently accurate in a small neighborhood of the current iterate $x^k$. The goal of this paper is to relax this requirement and allow the use of random local models which are accurate only with certain probability (conditioned on the past). In that case, note that the models need to be re-drawn after each iteration, whether successful or not.

Note that our general setting includes the cases when the model (the derivative information, for example) is always accurate, but the step $s^k$ is computed approximately, in a probabilistic manner. For example, $s^k$ can be an approximation of $-(b^k)^{-1}g^k$. It is easy to see how randomness in $s^k$ calculation can be viewed as the randomness in the model, by considering that instead of the accurate model

$$f(x^k) + (x - x^k)^T g^k + \frac{1}{2}(x - x^k)^\top b^k (x - x^k)$$

we use an approximate model

$$m_k(x) = f(x^k) - (x - x^k)^T b^k s^k + \frac{1}{2}(x - x^k)^\top b^k (x - x^k).$$

Hence, as long as the accuracy requirements are carried over accordingly the approximate random models subsume the case of approximate random step computations. The next section makes precise our requirements on the probabilistic models.

## 2.2 Generic probabilistic models

We will now introduce the key probabilistic ingredients of our scheme. In particular we assume that our models $m_k$ are random and that they satisfy some notion of good quality with some probability $p$. We will consider random models $M_k$, and then use the notation $m_k = M_k(\omega_k)$ for their realizations. The randomness of the models will imply the randomness of the points $x^k$, the step length parameter $\alpha_k$, the computed steps $s^k$ and other quantities produced by the algorithm. Thus, in our paper, these random variables will be denoted by $X^k$, $\mathcal{A}_k$, $S^k$ and so on, respectively, while $x^k = X^k(\omega_k)$, $\alpha_k = \mathcal{A}_k(\omega_k)$, $s^k = S^k(\omega_k)$, etc, denote their realizations (we will omit the $\omega_k$ in the notation for brevity).

For each specific optimization method, we will define a notion of sufficiently accurate models. The desired accuracy of the model depends on the current iterate $x^k$, step parameter $\alpha_k$ and, possibly, the step $s^k(\alpha_k)$. This notion involves model properties which make sufficient decrease in $f$ achievable by the step $s^k(\alpha_k)$. Specific conditions on the models will be stated for each algorithm in the respective sections and how these conditions may be achieved will be discussed in Section 5.

**Definition 2.1 [sufficiently accurate models; true and false iterations]** *We say that a sequence of random models $\{M_k\}$ is $(p)$-probabilistically "sufficiently accurate" for a corresponding sequence $\{\mathcal{A}_k, X^k\}$, if the following indicator random variable*

$$I_k \;=\; \mathbb{1}\{M_k \text{ is a sufficiently accurate model of } f \text{ for the given } X^k \text{and } \mathcal{A}_k\}$$

*satisfy the following submartingale-like condition*

$$P(I_k = 1 \,|\, \mathcal{F}_{k-1}^M) \;\geq\; p, \tag{1}$$

*where $\mathcal{F}_{k-1}^M = \sigma(M_0, \ldots, M_{k-1})$ is the $\sigma$-algebra generated by $M_0, \ldots, M_{k-1}$ - in other words, the history of the algorithm up to iteration $k$.*

*We say that iteration $k$ is a **true** iteration if the event $I_k = 1$ occurs. Otherwise the iteration is called **false**.*

Note that $M_k$ is a random model that, given the past history, encompasses all the randomness of iteration $k$ of our algorithm. The iterates $X^k$ and the step length parameter $\mathcal{A}_k$ are random variables defined over the $\sigma$-algebra generated by $M_0, \ldots, M_{k-1}$. Each $M_k$ depends on $X^k$ and $\mathcal{A}_k$ and hence on $M_0, \ldots, M_{k-1}$. Definition 2.1 serves to enforce the following property: even though the accuracy of $M_k$ may be dependent on the history, $(M_1, \ldots, M_{k-1})$, via its dependence on $X^k$ and $\mathcal{A}_k$, it is sufficiently good with probability at least $p$, regardless of that history. This condition is more reasonable than complete independence of $M^k$ from the past, which is difficult to ensure. It is important to note that, from this assumption, it follows that whether or not the step is deemed successful and the iterate $x^k$ is updated, our scheme always updates the model $m_k$, unless $m_k$ is somehow known to be sufficiently accurate for $x^{k+1} = x^k$ and $\alpha_{k+1}$. We will discuss this in more detail in Section 5.

When Algorithm 2.1 is based on probabilistic models (and all its specific variants under consideration), it results in a discrete time stochastic process. This stochastic process encompasses random elements such $\mathcal{A}_k$, $X^k$, $S^k$, which are directly computed by the algorithm, but also some quantities that can be derived as functions of $\mathcal{A}_k$, $X^k$, $S^k$, such as $f(X^k)$, $\|\nabla f(X^k)\|$ and a quantity $F_k$, which we will use to denote some measure of progress towards optimality. Each realization of the sequence of random models results in a realization of the algorithm, which in turn produces the corresponding sequences $\{\alpha_k\}$, $\{x^k\}$, $\{s^k\}$, $\{f(x^k)\}$, $\{\|\nabla f(x^k)\|\}$ and $\{f_k\}$[1]. We will analyze the stochastic processes restricting our attention to some of the random quantities that belong to this process and will ignore the rest, for the brevity of the presentation. Hence when we say that Algorithm 2.1 generates the stochastic process $\{X^k, \mathcal{A}_k\}$, this means we want to focus on the properties of these random variables, but keeping in mind that there are other random quantities in this stochastic process.

We will derive complexity bounds for each algorithm in the following sense. We will define the accuracy goal that we aim to reach and then we will bound the expected number of steps that the algorithm takes until this goal is achieved. The analyses will follow common steps, and the main ingredients are described below. We then apply these steps to each case under consideration.

## 2.3 Elements of global convergence rate analysis

First we recall a standard notion from stochastic processes.

---

[1] Note that throughout, $f(x^k) \neq f_k$, since $f_k$ is a related measure of progress towards optimality.

**Hitting time.** For a given discrete time stochastic process, $Z_t$, recall the concept of a *hitting time* for an event $\{Z_t \in S\}$. This is a random variable, defined as $T_S = \min\{t : Z_t \in S\}$ - the first time the event $\{Z_t \in S\}$ occurs. In our context, set $S$ will either be a set of real numbers larger than some given value, or smaller than some other given value.

**Number of iterations $N_\epsilon$ to reach $\epsilon$ accuracy.** Given a level of accuracy $\epsilon$, we aim to derive a bound on the expected number of iterations $\mathbb{E}(N_\epsilon)$ which occur in the algorithm until the given accuracy level is reached. The number of iterations $N_\epsilon$ is a random variable, which can be defined as a hitting time of some stochastic process, dependent on the case under analysis. In particular,

- If $f(x)$ is not known to be convex, then $N_\epsilon$ is the hitting time for $\{\|\nabla f(X_k)\| \leq \epsilon\}$, namely, the number of steps the algorithm takes until $\|\nabla f(X^k)\| \leq \epsilon$ occurs for the first time.

- If $f(x)$ is convex or strongly convex then $N_\epsilon$ is the hitting time for $\{f(X^k) - f_* \leq \epsilon\}$, namely, the number of steps the algorithm takes until $f(X^k) - f_* \leq \epsilon$ occurs for the first time, where $f_* = f(x^*)$ with $x^*$, a global minimizer of $f$.

We will bound $\mathbb{E}(N_\epsilon)$ by observing that for all $k < N_\epsilon$ the stochastic process induced by Algorithm 2.1 behaves in a certain way. To formalize this, we need to define the following random variable and its upper bound.

**Measure of progress towards optimality, $F_k$.** This measure is defined by the total function decrease or by the distance to the optimum. In particular,

- If $f(x)$ is not known to be convex, then $F_k = f(X^0) - f(X^k)$.

- If $f(x)$ is convex, then $F_k = 1/(f(X^k) - f_*)$.

- If $f(x)$ is strongly convex, then $F_k = \log(1/(f(X^k) - f_*))$.

**Upper bound $F_\epsilon$ on $F_k$.** From the algorithm construction, $F_k$ defined above is always non-decreasing and there exists a deterministic upper bound $F_\epsilon$ in each case, defined as follows.

- If $f(x)$ is not known to be convex, then $F_\epsilon = f(X^0) - f_*$, where $f_*$ is a global lower bound on $f$.

- If $f(x)$ is convex, then $F_\epsilon = 1/\epsilon$.

- If $f(x)$ is strongly convex, then $F_\epsilon = \log(1/\epsilon)$.

We observe that $F_k$ is a nondecreasing process and $F_\epsilon$ is the largest possible value that $F_k$ can achieve.

Our analysis will be based on the following observations, which are borrowed from the global rate analysis of the deterministic methods [15].

- **Guaranteed amount of increase in $f_k$.** For all $k < N_\epsilon$ (i.e., until the desired accuracy has been reached), if the $k$th iteration is true and successful, then $f_k$ is increased by an amount proportional to $\alpha_k$.

7

- **Guaranteed threshold for $\alpha_k$.** There exists a constant, which we will call $C$, such that if $\alpha_k \leq C$ and the $k$th iteration is true, then the $k$th iteration is also successful, and hence $\alpha_{k+1} = \gamma^{-1}\alpha_k$. This constant $C$ depends on the algorithm and Lipschitz constants of $f$.

- **Bound on the number of iterations.** If all iterations were true, then by the above observations, $\alpha_k \geq \gamma C$ and, hence, $f_k$ increases by at least a constant for all $k$. From this a bound on the number of iterations, knowing that $f^k$ cannot exceed $F_\epsilon$.

In our case not all iterations are true, however, under the assumption that they "tend" to be true, as we will show, when $\mathcal{A}_k \leq C$, then iterations "tend" to be successful, $\mathcal{A}_k$ "tends" to stay near the value $C$ and the values $F_k$ "tend" to increase by a constant. The analysis is then performed via a study of stochastic processes, which we describe in detail next.

## 2.4 Analysis of the stochastic processes

Let us consider the stochastic process $\{\mathcal{A}_k, F_k\}$ generated by Algorithm 2.1 using random, $p$-probabilistically sufficiently accurate models $M_k$, with $F_k$ defined above. Under the assumption that the sequence of models $M_k$ are $p$-probabilistically sufficiently accurate, each iteration is true with probability at least $p$, conditioned on the past.

We assume now (and we show later for each specific case) that $\{\mathcal{A}_k, F_k\}$ obeys the following rules for all $k < N_\epsilon$.

**Assumption 2.1** *There exist a constant $C > 0$ and a nondecreasing function $h(\alpha)$, $\alpha \in \mathbb{R}$, which satisfies $h(\alpha) > 0$ for any $\alpha > 0$, such that for any realization of Algorithm 2.1 the following hold for all $k < N_\epsilon$:*

- *(i) If iteration $k$ is true (i.e. $I_k = 1$) and successful, then $f_{k+1} \geq f_k + h(\alpha_k)$.*

- *(ii) If $\alpha_k \leq C$ and iteration $k$ is true then iteration $k$ is also successful, which implies $\alpha_{k+1} = \gamma^{-1}\alpha_k$.*

- *(iii) $f_{k+1} \geq f_k$ for all $k$.*

For future use let us state an auxiliary lemma.

**Lemma 2.1** *Let $N_\epsilon$ be the hitting time as defined on page 7. For all $k < N_\epsilon$, let $I_k$ be the sequence of random variables in Definition 2.1 so that (1) holds. Let $W_k$ be a nonnegative stochastic process such that $\sigma(W_k) \subset \mathcal{F}_{k-1}^M$, for any $k \geq 0$. Then*

$$\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} W_k I_k\right) \geq p\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} W_k\right).$$

*Similarly,*

$$\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} W_k(1 - I_k)\right) \leq (1-p)\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} W_k\right).$$

*Proof.* The proof is a simple consequence of properties of expectations, see for example, [22, property H*, page 216],

$$\mathbb{E}(I_k \,|\, W_k) = \mathbb{E}(\mathbb{E}(I_k \,|\, \mathcal{F}_{k-1}^M) \,|\, W_k) \geq \mathbb{E}(p \,|\, W_k) \geq p,$$

where we also used that $\sigma(W_k) \subset \mathcal{F}_{k-1}^M$. Hence by the law of total expectation, we have $\mathbb{E}(W_k I_k) = \mathbb{E}(W_k \mathbb{E}(I_k | W_k)) \geq p\mathbb{E}(W_k)$. Similarly, we can derive $\mathbb{E}(\mathbb{1}\{k < N_\epsilon\} W_k I_k) \geq p\mathbb{E}(\mathbb{1}\{k < N_\epsilon\} W_k)$, because $\mathbb{1}\{k < N_\epsilon\}$ is also determined by $\mathcal{F}_{k-1}^M$. Finally,

$$\mathbb{E}\left( \sum_{k=0}^{N_\epsilon - 1} W_k I_k \right) = \mathbb{E}\left( \sum_{k=0}^{\infty} \mathbb{1}\{k < N_\epsilon\} W_k I_k \right) \geq p\mathbb{E}\left( \sum_{k=0}^{\infty} \mathbb{1}\{k < N_\epsilon\} W_k \right) = p\mathbb{E}\left( \sum_{k=0}^{N_\epsilon - 1} W_k \right).$$

The second inequality is proved analogously. $\qquad\square$

Let us now define two indicator random variables, in addition to $I_k$ defined earlier,

$$\Lambda_k = \mathbb{1}\{\mathcal{A}_k > C\},$$

and

$$\Theta_k = \mathbb{1}\{\text{Iteration } k \text{ is successful i.e., } \mathcal{A}_{k+1} = \gamma^{-1}\mathcal{A}_k\}.$$

Note that $\sigma(\Lambda_k) \subset \mathcal{F}_{k-1}^M$ and $\sigma(\Theta_k) \subset \mathcal{F}_k^M$, that is the random variable $\Lambda_k$ is fully determined by the first $k-1$ steps of the algorithm, while $\Theta_k$ is fully determined by the first $k$ steps. We will use $\lambda_k$, $i_k$ and $\theta_k$ to denote realizations of $\Lambda_k$, $I_k$ and $\Theta_k$, respectively.

These indicators will help us define our algorithm more rigorously as a stochastic process. Without loss of generality, we assume that $C = \gamma^c \alpha_0 < \gamma \alpha_{\max}$ for some positive integer $c$. In other words, $C$ is the largest value that the step size $\mathcal{A}_k$ actually achieves for which part $(ii)$ of Assumption 2.1 holds. The condition $C < \gamma \alpha_{\max}$ is a simple technical condition, which is not necessary, but which simplifies the presentation later in this section. Under Assumption 2.1, recalling the update rules for $\alpha_k$ in Algorithm 2.1 and the assumption that true iterations occur with probability at least $p$, we can write the stochastic process $\{\mathcal{A}_k, F_k\}$ as obeying the expressions below:

$$\mathcal{A}_{k+1} = \begin{cases} \gamma^{-1}\mathcal{A}_k & \text{if } I_k = 1 \text{ and } \Lambda_k = 0, \\ \gamma\mathcal{A}_k & \text{if } I_k = 0 \text{ and } \Lambda_k = 0, \\ \min\{\alpha_{\max}, \gamma^{-1}\mathcal{A}_k\} & \text{if } \Theta_k = 1 \text{ and } \Lambda_k = 1, \\ \gamma\mathcal{A}_k & \text{if } \Theta_k = 0 \text{ and } \Lambda_k = 1, \end{cases} \qquad (2)$$

$$F_{k+1} \geq \begin{cases} F_k + h(\mathcal{A}_k) & \text{if } I_k = 1 \text{ and } \Lambda_k = 0, \\ F_k & \text{if } I_k = 0 \text{ and } \Lambda_k = 0, \\ F_k + h(\mathcal{A}_k) & \Theta_k I_k = 1 \text{ and } \Lambda_k = 1, \\ F_k & \Theta_k I_k = 0 \text{ and } \Lambda_k = 1. \end{cases} \qquad (3)$$

We conclude that, when $\mathcal{A}_k \leq C$, a successful iteration happens with probability at least $p$, and in that case $\mathcal{A}_{k+1} = \gamma^{-1}\mathcal{A}_k$, and that an unsuccessful iteration happens with probability at most $1-p$, in which case $\mathcal{A}_{k+1} = \gamma\mathcal{A}_k$. Note that there is no known probability bound for the different outcomes when $\mathcal{A}_k > C$. However, we know that $I_k = 1$ with probability at least $p$ and if, in addition, iteration $k$ happens to be successful, then $F_k$ is increased by at least $h(\mathcal{A}_k)$.

In summary, from the above discussion, we have

*for all $k < N_\epsilon$, Algorithm 2.1 under Assumption 2.1 yields the stochastic process $\{\mathcal{A}_k, F_k\}$ in (2) and (3).*

9

## 2.5 Bounding the number of steps for which $\alpha_k \leq C$

In this subsection we derive a bound on $\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1}(1-\Lambda_k)\right)$. The bound for $\mathbb{E}(\sum_{k=0}^{N_\epsilon-1}\Lambda_k)$ will be derived in the next section.

The following simple result holds for every realization of the algorithm and stochastic process $\{\Lambda_k, I_k, \Theta_k\}$.

**Lemma 2.2** *For any $l \in \{0, \ldots, N_\epsilon - 1\}$ and for all realizations of Algorithm 2.1, we have*

$$\sum_{k=0}^{l}(1-\Lambda_k)\Theta_k \leq \frac{1}{2}(l+1).$$

*Proof.* By the definition of $\Lambda_k$ and $\Theta_k$ we know that when $(1-\Lambda_k)\Theta_k = 1$ then we have a successful iteration and $\mathcal{A}_k \leq C$. In this case $\mathcal{A}_{k+1} = \gamma^{-1}\mathcal{A}_k$. It follows that amongst all iterations, at most half can be successful and have $\mathcal{A}_k \leq C$, because for each such iteration, when $\mathcal{A}_k$ gets increased by a factor of $\gamma^{-1}$, there has to be at least one iteration when $\mathcal{A}_k$ is decreased by the same factor, since $\mathcal{A}_0 \geq C$. $\qquad\square$

Using this we derive the bound.

**Lemma 2.3**
$$\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1}(1-\Lambda_k)\right) \leq \frac{1}{2p}\mathbb{E}(N_\epsilon)$$

*Proof.* By Lemma 2.1 applied to $W_k = 1 - \Lambda_k$ we have

$$\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1}(1-\Lambda_k)I_k\right) \geq p\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1}(1-\Lambda_k)\right). \tag{4}$$

From the fact that all true iterations are successful when $\alpha_k \leq C$,

$$\sum_{k=0}^{N_\epsilon-1}(1-\Lambda_k)I_k \leq \sum_{k=0}^{N_\epsilon-1}(1-\Lambda_k)\Theta_k. \tag{5}$$

Finally, from Lemma 2.2

$$\sum_{k=0}^{N_\epsilon-1}(1-\Lambda_k)I_k \leq \frac{1}{2}N_\epsilon. \tag{6}$$

Taking expectations in (5) and (6) and combining with (4), we obtain the result of the lemma. $\qquad\square$

## 2.6 Bounding the expected number of steps for which $\alpha_k > C$

Let us now consider the bound on $\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1}\Lambda_k\right)$. We introduce the additional notation $\bar{\Lambda}_k = \mathbb{1}\{\mathcal{A}_k > C\} + \mathbb{1}\{\mathcal{A}_k = C\}$. In other words $\bar{\Lambda}_k = 1$ when either $\Lambda_k = 1$ or $\mathcal{A}_k = C$. We now define:

- $N_1 = \sum_{k=0}^{N_\epsilon - 1} \bar{\Lambda}_k (1 - I_k) \Theta_k$, which is the number of false successful iterations, when $\mathcal{A}_k \geq C$.

- $M_1 = \sum_{k=0}^{N_\epsilon - 1} \bar{\Lambda}_k (1 - I_k)$, which is the number of false iterations, when $\mathcal{A}_k \geq C$.

- $N_2 = \sum_{k=0}^{N_\epsilon - 1} \bar{\Lambda}_k I_k \Theta_k$, which is the number of true successful iterations, when $\mathcal{A}_k \geq C$.

- $M_2 = \sum_{k=0}^{N_\epsilon - 1} \bar{\Lambda}_k I_k$, which is the number of true iterations, when $\mathcal{A}_k \geq C$.

- $N_3 = \sum_{k=0}^{N_\epsilon - 1} \Lambda_k I_k (1 - \Theta_k)$, which is the number of true unsuccessful iterations, when $\mathcal{A}_k > C$.

- $M_3 = \sum_{k=0}^{N_\epsilon - 1} \Lambda_k (1 - \Theta_k)$, which is the number of unsuccessful iterations, when $\mathcal{A}_k > C$.

Since $\mathbb{E}\left(\sum_{k=0}^{N_\epsilon - 1} \Lambda_k\right) = \mathbb{E}\left(\sum_{k=0}^{N_\epsilon - 1} \Lambda_k (1 - I_k)\right) + \mathbb{E}\left(\sum_{k=0}^{N_\epsilon - 1} \Lambda_k I_k\right) \leq \mathbb{E}(M_1) + \mathbb{E}(M_2)$, our goal is to bound $\mathbb{E}(M_1) + \mathbb{E}(M_2)$.

Our next observation is simple but central in our analysis. It reflects the fact that the gain in $F_k$ is bounded from above by $F_\epsilon$ and when $\mathcal{A}_k \geq C$ this gain is bounded from below as well, hence allowing us to bound the total number of true successful iterations when $\mathcal{A}_k \geq C$. The following two lemmas holds for every realization.

**Lemma 2.4** *For any $l \in \{0, \ldots, N_\epsilon - 1\}$ and for all realizations of Algorithm 2.1, we have*

$$\sum_{k=0}^{l} \bar{\Lambda}_k I_k \Theta_k \leq \frac{F_\epsilon}{h(C)},$$

*and so*

$$N_2 \leq \frac{F_\epsilon}{h(C)}. \tag{7}$$

*Proof.* Consider any $k$ for which $\Lambda_k I_k \Theta_k = 1$. From Assumption 2.1 we know that whenever an iteration is true and successful then $F_k$ get increased by at least $h(\mathcal{A}_k) \geq h(C)$, since $\mathcal{A}_k \geq C$ and $h$ is nondecreasing. We also know that on other iterations $F_k$ does not decrease. The bound $F_k \leq F_\epsilon$ trivially gives us the desired result. $\square$

Another key observation is that

$$M_2 \leq N_2 + N_3 \leq N_2 + M_3, \tag{8}$$

where the first inequality follows from the fact that for all $k < N_\epsilon$ and for all realizations, $(\bar{\Lambda}_k - \Lambda_k) I_k (1 - \Theta_k) = 0$, in other words there are no true unsuccessful iterations when $\mathcal{A}_k = C$.

**Lemma 2.5** *For any $l \in \{0, \ldots, N_\epsilon - 1\}$ and for all realizations of Algorithm 2.1, we have*

$$\sum_{k=0}^{l} \Lambda_k (1 - \Theta_k) \leq \sum_{k=0}^{l} \bar{\Lambda}_k \Theta_k + \log_\gamma \left(\frac{C}{\alpha_0}\right)$$

*Proof.* $\mathcal{A}_k$ is increased on successful iterations and decreased on unsuccessful ones. Hence the total number of steps when $\mathcal{A}_k > C$ and $\mathcal{A}_k$ is decreased, is bounded by the total number of steps when $\mathcal{A}_k \geq C$ is increased plus the number of steps it is required to reduce $\mathcal{A}_k$ from its initial value $\alpha_0$ to $C$. □

From Lemma 2.5 applied to $l = N_\epsilon - 1$, we can deduce that

$$M_3 \leq N_1 + N_2 + \log_\gamma(C/\alpha_0). \tag{9}$$

We also have the following lemma.

**Lemma 2.6**

$$\mathbb{E}(M_1) \leq \frac{1-p}{p}\mathbb{E}(M_2). \tag{10}$$

*Proof.* By applying both inequalities in Lemma 2.1 with $W_k = \bar{\Lambda}_k$, we obtain

$$\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} \bar{\Lambda}_k I_k\right) \geq p\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} \bar{\Lambda}_k\right)$$

and

$$\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} \bar{\Lambda}_k(1 - I_k)\right) \leq (1-p)\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} \bar{\Lambda}_k\right)$$

which gives us

$$\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} \bar{\Lambda}_k(I_k - 1)\right) \leq \frac{1-p}{p}\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} \bar{\Lambda}_k I_k\right).$$

□

**Lemma 2.7** *Under the condition that $p > 1/2$, we have*

$$\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1} \Lambda_k\right) \leq \frac{2F_\epsilon}{h(C)(2p-1)} + \frac{\log_\gamma(C/\alpha_0)}{2p-1}.$$

*Proof.* Recall that $\mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1}\Lambda_k\right) = \mathbb{E}(M_1 + M_2)$. Using (8) and (10) it follows that

$$\mathbb{E}(N_1) \leq \mathbb{E}(M_1) \leq \frac{1-p}{p}\mathbb{E}(M_2) \leq \frac{1-p}{p}\mathbb{E}(N_2 + M_3) = \frac{1-p}{p}[\mathbb{E}(N_2) + \mathbb{E}(M_3)]. \tag{11}$$

Taking into account (9) and using the bound (7) on $N_2$ we have

$$\mathbb{E}(M_3) \leq \mathbb{E}(N_1) + \mathbb{E}(N_2) + \log_\gamma(C/\alpha_0) \leq \mathbb{E}(N_1) + F_\epsilon/h(C) + \log_\gamma(C/\alpha_0). \tag{12}$$

Plugging this into (11) and using the bound (7) on $N_2$ again, we obtain

$$\mathbb{E}(N_1) \leq \frac{1-p}{p}\left[\frac{F_\epsilon}{h(C)} + \mathbb{E}(N_1) + \frac{F_\epsilon}{h(C)} + \log_\gamma\left(\frac{C}{\alpha_0}\right)\right],$$

and, hence,

$$\frac{2p-1}{p}\mathbb{E}(N_1) \leq \frac{1-p}{p}\left[\frac{2F_\epsilon}{h(C)} + \log_\gamma\left(\frac{C}{\alpha_0}\right)\right].$$

This finally implies

$$\mathbb{E}(N_1) \leq \frac{1-p}{2p-1}\left[\frac{2F_\epsilon}{h(C)} + \log_\gamma\left(\frac{C}{\alpha_0}\right)\right]. \tag{13}$$

Now we can bound the expected total number of iterations when $\alpha_k > C$, using (7), (12) and (13) and adding the terms to obtain the result of the lemma, namely,

$$\mathbb{E}(M_1 + M_2) \leq \mathbb{E}(M_1 + M_3 + N_2) \leq \frac{1}{p}\mathbb{E}(M_3 + N_2) \leq \frac{1}{2p-1}\left(\frac{2F_\epsilon}{h(C)} + \log_\gamma\left(\frac{C}{\alpha_0}\right)\right).$$

$\square$

## 2.7 Final bound on the expected stopping time

We finally have the following theorem which trivially follows from Lemmas 2.3 and 2.7.

**Theorem 2.1** *Under the condition that $p > 1/2$, the hitting time $N_\epsilon$ is bounded in expectation as follows*

$$\mathbb{E}(N_\epsilon) \leq \frac{2p}{(2p-1)^2}\left(\frac{2F_\epsilon}{h(C)} + \log_\gamma\left(\frac{C}{\alpha_0}\right)\right).$$

*Proof.* Clearly

$$\mathbb{E}(N_\epsilon) = \mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1}\Lambda_k\right) + \mathbb{E}\left(\sum_{k=0}^{N_\epsilon-1}(1-\Lambda_k)\right)$$

and, hence, using Lemmas 2.3 and 2.7 we have

$$\mathbb{E}(N_\epsilon) \leq \frac{1}{2p}\mathbb{E}(N_\epsilon) + \frac{1}{2p-1}\left(\frac{2F_\epsilon}{h(C)} + \log_\gamma\left(\frac{C}{\alpha_0}\right)\right).$$

The result of the theorem easily follows. $\square$

**Summary of our complexity analysis framework.** We have considered a(ny) algorithm in the framework Algorithm 2.1 with probabilistically sufficiently accurate models as in Definition 2.1. We have developed a methodology to obtain (complexity) bounds on the number of iterations $N_\epsilon$ that such an algorithm takes to reach desired accuracy. It is important to note that, while we simply provide the bound on $\mathbb{E}(N_\epsilon)$ it is easy to extend the analysis of the same stochastic processes to provide bounds on $P\{N_\epsilon > K\}$, for any $K$ larger than the bound on $\mathbb{E}(N_\epsilon)$, in particular it can be shown that $P\{N_\epsilon > K\}$ decays exponentially with $K$.

While in our analysis we assumed that the constant $\gamma$ by which we decrease and increase $\alpha_k$ is the same, our analysis can be quite easily extended to the case when the constants for increase and decrease are different, say $\gamma_{inc}$ and $\gamma_{dec}$. In this case the threshold on the probability $p$ may no longer be $1/2$ but will be larger if $\gamma_{inc}/\gamma_{dec} < 1$ and smaller, otherwise. Some of the constants in the upper bound on $\mathbb{E}(N_\epsilon)$ with change accordingly.

Our approach is valid provided that all of the conditions in Assumption 2.1 hold. Next we show that all these conditions are satisfied by steepest-descent linesearch methods in the nonconvex, convex and strongly convex case; by general linesearch methods in the nonconvex case; by cubic regularization methods (ARC) for nonconvex objectives. In particular, we will specify what we mean by a probablistically sufficiently accurate first-order and second-order model in the case of linesearch and cubic regularization methods, respectively.

# 3   The line-search algorithm

We will now apply the generic analysis outlined in the previous section to the case of the following simple probabilistic line-search algorithm.

**Algorithm 3.1 A line-search algorithm with random models**

**Initialization**
   Choose constants $\gamma \in (0,1)$, $\theta \in (0,1)$ and $\alpha_{\max} > 0$. Pick initial $x^0$ and $\alpha_0 < \alpha_{\max}$.
   Repeat for $k = 0, 1, \ldots$

1. **Compute a model and a step**
   Compute a random model $m_k$ and use it to generate a direction $g^k$.
   Set the step $s^k = -\alpha_k g^k$.

2. **Check sufficient decrease**
   Check if
   $$f(x^k - \alpha_k g^k) \le f(x^k) - \alpha_k \theta \|g^k\|^2. \tag{14}$$

3. **Successful step**
   If (14) holds, then $x^{k+1} := x^k - \alpha_k g^k$ and $\alpha_{k+1} = \min\{\alpha_{\max}, \gamma^{-1}\alpha_k\}$. Let $k := k + 1$.

4. **Unsuccessful step**
   Otherwise, $x^{k+1} := x^k$, set $\alpha_{k+1} = \gamma \alpha_k$. Let $k := k + 1$.

For the linesearch algorithm, the key ingredient is a search direction selection on each iteration. In our case we assume that the search direction is random and satisfies some accuracy requirement that we discuss below. The choice of model in this algorithm is a simple linear model $m_k(x)$, which gives rise to the search direction $g^k$, specifically, $m_k(x) = f(x^k) + (x - x^k)^T g^k$. We will consider more general models in the next section, Section 3.2.

Recall Definition 2.1. Here we describe the specific requirement we apply to the models in the case of line search.

**Definition 3.1** *We say that a sequence of random models and corresponding directions $\{M_k, G_k\}$ is $(p)$-probabilistically "sufficiently accurate" for Algorithm 3.1 for a corresponding sequence $\{\mathcal{A}_k, X^k\}$, if there exists a constant $\kappa > 0$, such that the indicator variables*

$$I_k \; = \; \mathbb{1}\{\|G^k - \nabla f(X^k)\| \le \kappa \mathcal{A}_k \|G^k\|\}$$

*satisfy the following submartingale-like condition*

$$P(I_k = 1 | \mathcal{F}_{k-1}^M) \; \ge \; p,$$

*where $\mathcal{F}_{k-1}^M = \sigma(M_0, \ldots, M_{k-1})$ is the $\sigma$-algebra generated by $M_0, \ldots, M_{k-1}$.*

As before, each iteration for which $I_k = 1$ holds is called a true iteration. It follows that for every realization of the algorithm, on all true iterations, we have

$$\|g^k - \nabla f(x^k)\| \leq \kappa \alpha_k \|g^k\|, \tag{15}$$

which implies, using $\alpha_k \leq \alpha_{\max}$ and the triangle inequality, that

$$\|g^k\| \geq \frac{\|\nabla f(x^k)\|}{1 + \kappa \alpha_{\max}}. \tag{16}$$

For the remainder of the analysis of Algorithm 3.1, we make the following assumption.

**Assumption 3.1** *The sequence of random models and corresponding directions $\{M_k, G_k\}$, generated in Algorithm 3.1, is $(p)$-probabilistically "sufficiently accurate" for the corresponding random sequence $\{\mathcal{A}_k, X^k\}$, with $p > 1/2$.*

We also make a standard assumption on the smoothness of $f(x)$ for the remainder of the paper.

**Assumption 3.2** $f \in \mathcal{C}^1(\mathbb{R}^n)$, *is globally bounded below by $f_*$, and has globally Lipschitz continuous gradient $\nabla f$, namely,*

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \quad \text{for all } x, y \in \mathbb{R}^n \text{ and some } L > 0. \tag{17}$$

## 3.1 The nonconvex case, steepest descent

As mentioned before, our goal in the nonconvex case is to compute a bound on the expected number of iterations $k$ that Algorithm 3.1 requires to obtain an iterate $x^k$ for which $\|\nabla f(x^k)\| \leq \epsilon$. We will now compute the specific quantities and expressions defined in Sections 2.3 and 2.4, that allow us to apply the analysis of our general framework to the specific case of Algorithm 3.1 for nonconvex functions.

Let $N_\epsilon$ denote, as before, the number of iterations that are taken until $\|\nabla f(X^k)\| \leq \epsilon$ occurs (which is a random variable). Let us consider the stochastic process $\{\mathcal{A}_k, F_k\}$ with $F_k = f(x^0) - f(X^k)$ and let $F_\epsilon = f(x^0) - f_*$. Then $F_k \leq F_\epsilon$, for all $k$.

Next we show that Assumption 2.1 is verified. First we derive an expression for the constant $C$, related to the size of the stepsize $\alpha_k$.

**Lemma 3.1** *Let Assumption 3.2 hold. For every realization of Algorithm 3.1, if iteration $k$ is true (i.e. $I_k = 1$), and if*

$$\alpha_k \leq C = \frac{1 - \theta}{0.5L + \kappa}, \tag{18}$$

*then (14) holds. In other words, when (18) holds, any true iteration is also a successful one.*

*Proof.* Condition (17) implies the following overestimation property for all $x$ and $s$ in $\mathbb{R}^n$,

$$f(x + s) \leq f(x) + s^T \nabla f(x) + \frac{L}{2}\|s\|^2,$$

which implies

$$f(x^k - \alpha_k g^k) \leq f(x^k) - \alpha_k (g^k)^T \nabla f(x^k) + \tfrac{L}{2}\alpha_k^2 \|g^k\|^2.$$

Applying the Cauchy-Schwarz inequality and (15) we have

$$
\begin{aligned}
f(x^k - \alpha_k g^k) &\leq f(x^k) - \alpha_k (g^k)^T [\nabla f(x^k) - g^k] - \alpha_k \|g^k\|^2 \left[1 - \tfrac{L}{2}\alpha_k\right] \\
&\leq f(x^k) + \alpha_k \|g^k\| \cdot \|\nabla f(x^k) - g^k\| - \alpha_k \|g^k\|^2 \left[1 - \tfrac{L}{2}\alpha_k\right] \\
&\leq f(x^k) - \alpha_k \|g^k\|^2 \left[1 - \left(\kappa + \tfrac{L}{2}\right)\alpha_k\right].
\end{aligned}
$$

It follows that (14) holds whenever $f(x^k) - \alpha_k \|g^k\|^2 [1 - (\kappa + 0.5L)\alpha_k] \leq f(x^k) - \alpha_k \theta \|g^k\|^2$ which is equivalent to (18). $\qquad\square$

From Lemma 3.1, and from (14) and (16), for any realization of Algorithm 3.1 which gives us the specific sequence $\{\alpha_k, f_k\}$, the following hold.

- If $k$ is a true and successful iteration, then

$$
f_{k+1} \geq f_k + \frac{\theta \|\nabla f(x^k)\|^2 \alpha_k}{(1 + \kappa \alpha_{\max})^2}
$$

  and

$$
\alpha_{k+1} = \gamma^{-1} \alpha_k.
$$

- If $\alpha_k \leq C$, where $C$ is defined in (18), and iteration $k$ is true, then it is also successful.

Hence, Assumption 2.1 holds and the process $\{\mathcal{A}_k, F_k\}$ behaves exactly as our generic process (2)-(3) in Section 2.4, with $C$ defined in (18) and the specific choice of $h(\mathcal{A}_k) = \frac{\theta \epsilon^2 \mathcal{A}_k}{(1 + \kappa \alpha_{\max})^2}$.

Finally, we use Theorem 2.1 and substituting the expressions for $C$, $h(C)$ and $F_\epsilon$ into the bound on $\mathbb{E}(N_\epsilon)$ we obtain the following complexity result.

**Theorem 3.1** *Let Assumptions 3.1 and 3.2 hold. Then the expected number of iterations that Algorithm 3.1 takes until $\|\nabla f(X^k)\| \leq \epsilon$ occurs is bounded as follows*

$$
\mathbb{E}(N_\epsilon) \leq \frac{2p}{(2p-1)^2} \left[\frac{M}{\epsilon^2} + \log_\gamma \left(\frac{1-\theta}{\alpha_0(0.5L + \kappa)}\right)\right],
$$

*where $M = \frac{(f(x^0) - f_*)(1 + \kappa \alpha_{\max})^2(0.5L + \kappa)}{\theta(1 - \theta)}$ is a constant independent of $p$ and $\epsilon$.*

**Remark 3.1** *We note that the dependency of the expected number of iterations on $\epsilon$ is of the order $1/\epsilon^2$, as expected from a line-search method applied to a smooth nonconvex problem. The dependency on $p$ is rather intuitive as well: if $p = 1$, then the deterministic complexity is recovered, while as $p$ approaches $1/2$, the expected number of iterations goes to infinity, since the models/directions are arbitrarily bad as often as they are good.*

## 3.2 The nonconvex case, general descent

In this subsection, we explain how the above analysis of the line-search method extends from the nonconvex steepest descent case to a general nonconvex descent case.

In particular, we consider that in Algorithm 3.1, $s^k = \alpha_k d^k$ (instead of $-\alpha_k g^k$), where $d^k$ is any direction that satisfies the following standard conditions.

- There exists a constant $\beta > 0$, such that

$$\frac{(d^k)^T g^k}{\|d^k\| \cdot \|g^k\|} \leq -\beta, \quad \forall k. \tag{19}$$

- There exist constants $\kappa_1, \kappa_2 > 0$, such that

$$\kappa_1 \|g^k\| \leq \|d^k\| \leq \kappa_2 \|g^k\|, \quad \forall k. \tag{20}$$

The sufficient decrease condition (14) is replaced by

$$f(x^k + \alpha_k d^k) \leq f(x^k) + \alpha_k \theta (d^k)^T g^k. \tag{21}$$

It is easy to show that a simple variant of Lemma 3.1 applies.

**Lemma 3.2** *Let Assumption 3.2 hold. Consider Algorithm 3.1 with $s^k = \alpha_k d^k$ and sufficient decrease condition (21). Assume that $d^k$ satisfies (19) and (20). Then, for every realization of the resulting algorithm, if iteration $k$ is true (i.e. $I_k$ holds), and if*

$$\alpha_k \leq C = \frac{\beta(1-\theta)}{0.5L\kappa_2 + \kappa}, \tag{22}$$

*then (21) holds. In other words, when (22) holds, any true iteration is also a successful one.*

*Proof.* The first displayed equation in the proof of Lemma 3.1 provides

$$f(x^k + \alpha_k d^k) \leq f(x^k) + \alpha_k (d^k)^T \nabla f(x^k) + \tfrac{L}{2}\alpha_k^2 \|d^k\|^2.$$

Applying the Cauchy-Schwarz inequality, (15) and the conditions (20) on $d^k$ we have

$$
\begin{aligned}
f(x^k + \alpha_k d^k) &\leq f(x^k) + \alpha_k (d^k)^T[\nabla f(x^k) - g^k] + \alpha_k (d^k)^T g^k + \tfrac{L}{2}\alpha_k^2 \|d^k\|^2 \\
&\leq f(x^k) + \alpha_k \|d^k\| \cdot \|\nabla f(x^k) - g^k\| + \alpha_k (d^k)^T g^k + \tfrac{L}{2}\alpha_k^2 \|d^k\|^2 \\
&\leq f(x^k) + \alpha_k^2 \kappa \|d^k\| \|g^k\| + \alpha_k (d^k)^T g^k + \tfrac{L}{2}\alpha_k^2 \kappa_2 \|d^k\| \|g^k\| \\
&= f(x^k) + \alpha_k (d^k)^T g^k + \alpha_k^2 \|d^k\| \|g^k\| \left(\kappa + \kappa_2 \tfrac{L}{2}\right).
\end{aligned}
$$

It follows that (21) holds whenever

$$\alpha_k (d^k)^T g^k + \alpha_k^2 \|d^k\| \|g^k\| \left(\kappa + \kappa_2 \frac{L}{2}\right) \leq \alpha_k \theta (d^k)^T g^k,$$

or equivalently, since $\alpha_k > 0$, whenever

$$\alpha_k \|d^k\| \|g^k\| \left(\kappa + \kappa_2 \frac{L}{2}\right) \leq -(1-\theta)(d^k)^T g^k.$$

Using (19), the latter displayed equation holds whenever $\alpha_k$ satisfies (22). $\qquad\square$

We conclude this extension to general descent directions by observing that if $k$ is a true and successful iteration, using the sufficient decrease condition (21), the conditions (19) and (20) on $d^k$ and (16), we obtain that

$$f_{k+1} \geq f_k + \frac{\theta \kappa_1 \beta \|\nabla f(x^k)\|^2 \alpha_k}{(1 + \kappa \alpha_{\max})^2}.$$

Hence, Assumption 2.1 holds for this case as well and the remainder of the analysis is exactly the same as for the steepest descent case.

## 3.3 The convex case

We now analyze the expected complexity of Algorithm 3.1 in the case when $f(x)$ is a convex function, that is when the following assumption holds.

**Assumption 3.3**   $f \in \mathcal{C}^1(\mathbb{R}^n)$ *is convex and has bounded level sets so that*

$$\|x - x^*\| \leq D \quad \text{for all } x \text{ with } f(x) \leq f(x^0), \tag{23}$$

*where $x^*$ is a global minimizer of $f$. Let $f^* = f(x^*)$.*

In this case, our goal is to bound the expectation of $N_\epsilon$ - the number of iterations taken by Algorithm 3.1 until

$$f(X^k) - f^* \leq \epsilon \tag{24}$$

occurs. We denote $f(X^k) - f^*$ by $\Delta_k^f$ and define $F_k = \frac{1}{\Delta_k^f}$. Clearly, $N_\epsilon$ is also the number of iterations taken until $F_k \geq \frac{1}{\epsilon} = F_\epsilon$ occurs.

Regarding Assumption 2.1, Lemma 3.1 provides the value for the constant $C$, namely, that whenever $\mathcal{A}_k \leq C$ with $C = \frac{1-\theta}{0.5L+\kappa}$, then every true iteration is also successful. We now show that on true and successful iterations, $F_k$ is increased by at least some function value $h(\mathcal{A}_k)$ for all $k < N_\epsilon$.

**Lemma 3.3** *Let Assumptions 3.2 and 3.3 hold. Consider any realization of Algorithm 3.1. For every iteration $k$ that is true and successful, we have*

$$f_{k+1} \geq f_k + \frac{\theta \alpha_k}{D^2(1 + \kappa \alpha_{\max})^2}. \tag{25}$$

*Proof.* Note that convexity of $f$ implies that for all $x$ and $y$,

$$f(x) - f(y) \geq \nabla f(y)^T (x - y),$$

and so by using $x = x^*$ and $y = x^k$, we have

$$-\Delta_k^f = f(x^*) - f(x^k) \geq \nabla f(x^k)^T (x^* - x^k) \geq -D\|\nabla f(x^k)\|,$$

where to obtain the last inequality, we used Cauchy-Schwarz inequality and (23). Thus when $k$ is a true iteration, (16) further provides

$$\frac{1}{D}\Delta_k^f \leq \|\nabla f(x^k)\| \leq (1 + \kappa \alpha_{\max})\|g^k\|.$$

When $k$ is also successful,

$$\Delta_k^f - \Delta_{k+1}^f = f(x^k) - f(x^{k+1}) \geq \theta \alpha_k \|g^k\|^2 \geq \frac{\theta \alpha_k}{D^2(1 + \kappa \alpha_{\max})^2}(\Delta_k^f)^2.$$

Dividing the above expression by $\Delta_k^f \Delta_{k+1}^f$, we have that on all true and successful iterations

$$\frac{1}{\Delta_{k+1}^f} - \frac{1}{\Delta_k^f} \geq \frac{\theta \alpha_k}{D^2(1 + \kappa \alpha_{\max})^2} \frac{\Delta_k^f}{\Delta_{k+1}^f} \geq \frac{\theta \alpha_k}{D^2(1 + \kappa \alpha_{\max})^2},$$

since $\Delta_k^f \geq \Delta_{k+1}^f$. Recalling the definition of $f_k$ completes the proof.   $\square$

Similarly to the nonconvex case, we conclude from Lemmas 3.1 and 3.3, that for any realization of Algorithm 3.1 the following have to happen.

- If $k$ is a true and successful iteration, then

$$f_{k+1} \geq f_k + \frac{\theta \alpha_k}{D^2(1 + \kappa \alpha_{\max})^2}$$

and

$$\alpha_{k+1} = \gamma^{-1} \alpha_k.$$

- If $\alpha_k \leq C$, where $C$ is defined in (18), and iteration $k$ is true, then it is also successful.

Hence, Assumption 2.1 holds and the process $\{\mathcal{A}_k, F_k\}$ behaves exactly as our generic process (2)-(3) in Section 2.4, with $C$ defined in (18) and the specific choice of $h(\mathcal{A}_k) = \frac{\theta \mathcal{A}_k}{D^2(1+\kappa\alpha_{\max})^2}$.

Theorem 2.1 can be immediately applied together with the above expressions for $C$, $h(C)$ and $F_\epsilon$, yielding the following complexity bound.

**Theorem 3.2** *Let Assumptions 3.1, 3.2 and 3.3 hold. Then the expected number of iterations that Algorithm 3.1 takes until $f(X^k) - f^* \leq \epsilon$ occurs is bounded by*

$$\mathbb{E}(N_\epsilon) \leq \frac{2p}{(2p-1)^2} \left[ \frac{M}{\epsilon} + \log_\gamma \left( \frac{1-\theta}{\alpha_0(0.5L + \kappa)} \right) \right],$$

*where $M = \frac{(1+\kappa\alpha_{\max})^2 D^2(0.5L+\kappa)}{\theta(1-\theta)}$ is a constant independent of $p$ and $\epsilon$.*

**Remark 3.2** *We again note the same dependence on $\epsilon$ in the complexity bound in Theorem 3.2 as in the deterministic convex case and on $p$, as in the nonconvex case.*

## 3.4 The strongly convex case

We now consider the case of strongly convex objective functions, hence the following assumption holds.

**Assumption 3.4** $f \in \mathcal{C}^1(\mathbb{R}^n)$ *is strongly convex, namely, for all $x$ and $y$ and some $\mu > 0$,*

$$f(x) \geq f(y) + \nabla f(y)^T (x - y) + \frac{\mu}{2} \|x - y\|^2.$$

Recall our notation $\Delta_k^f = f(X^k) - f^*$. Our goal here is again, as in the convex case, to bound the expectation on the number of iteration that occur until $\Delta_k^f \leq \epsilon$. In the strongly convex case, however, this bound is logarithmic in $\frac{1}{\epsilon}$, just as it is in the case of the deterministic algorithm.

**Lemma 3.4** *Let Assumption 3.4 hold. Consider any realization of Algorithm 3.1. For every iteration $k$ that is true and successful, we have*

$$f(x^k) - f(x^{k+1}) = \Delta_k^f - \Delta_{k+1}^f \geq \frac{2\mu\theta}{(1 + \kappa\alpha_{\max})^2} \alpha_k \Delta_k^f, \tag{26}$$

*or equivalently,*

$$\Delta_{k+1}^f \leq \left( 1 - \frac{2\mu\theta}{(1 + \kappa\alpha_{\max})^2} \alpha_k \right) \Delta_k^f. \tag{27}$$

*Proof.* Assumption 3.4 implies, for $x = x^k$ and $y = x^*$, that [see [15], Th 2.1.10]

$$\Delta_k^f \leq \frac{1}{2\mu} \|\nabla f(x^k)\|^2$$

or equivalently,

$$\sqrt{2\mu\Delta_k^f} \leq \|\nabla f(x^k)\| \leq (1 + \kappa\alpha_{\max})\|g^k\|,$$

where in the second inequality we used (16). The bound (26) now follows from the sufficient decrease condition (14). $\qquad\square$

Note that from (26) we have that if $\Delta_k^f > 0$ and $\alpha_k > (1 + \kappa\alpha_{\max})^2/(2\mu\theta)$ then the iteration is unsuccessful. Hence, for an iteration to be successful we must have $\alpha_k \leq (1 + \kappa\alpha_{\max})^2/(2\mu\theta)$. We also know that a true iteration is successful when $\alpha_k \leq C$, where $C$ defined in (18), assuming that $C \leq (1 + \kappa\alpha_{\max})^2/(2\mu\theta)$. To simplify the analysis we will simply assume that this inequality holds, by an appropriate choice of the parameters, which can done without loss of generality.

We now define $F_k = \log \frac{1}{\Delta_k^f}$ and $F_\epsilon = \log \frac{1}{\epsilon}$, and the hitting time $N_\epsilon$ is the number of iterations taken until $\Delta_k^f \leq \epsilon$.

As in the convex case, using Lemmas 3.1 and 3.4, we conclude that, for any realization of Algorithm 3.1, the following have to happen.

- If $k$ is a true and successful iteration, then

$$f_{k+1} \geq f_k - \log\left(1 - \frac{2\mu\theta}{(1 + \kappa\alpha_{\max})^2}\alpha_k\right).$$

  and

$$\alpha_{k+1} = \gamma^{-1}\alpha_k.$$

- If $\alpha_k \leq C$, where $C$ defined in (18), and iteration $k$ is true, then it is also successful.

Hence, again, Assumption 2.1 holds and the process $\{\mathcal{A}_k, F_k\}$ behaves exactly as our generic process (2)-(3) in Section 2.4, with $C$ defined in (18) and the specific choice of

$$h(\mathcal{A}_k) = -\log\left(1 - \frac{2\mu\theta}{(1 + \kappa\alpha_{\max})^2}\mathcal{A}_k\right).$$

By using the above expressions for $C$, $h(C)$ and $F_\epsilon$, again as in the convex case, we have the following complexity bound for the strongly convex case.

**Theorem 3.3** *Let Assumptions 3.1, 3.2 and 3.4 hold. Then the expected number of iterations that Algorithm 3.1 takes until $f(X^k) - f^* \leq \epsilon$ occurs is bounded by*

$$\mathbb{E}(N_\epsilon) \leq \frac{2p}{(2p-1)^2}\left[M\log\left(\frac{1}{\epsilon}\right) + \log_\gamma\left(\frac{1-\theta}{\alpha_0(0.5L + \kappa)}\right)\right],$$

*where $M = -\log\left(1 - \frac{2\mu\theta(1-\theta)}{(1+\kappa\alpha_{\max})^2(0.5L+\kappa)}\right)$ is a constant independent of $p$ and $\epsilon$.*

**Remark 3.3** *Again, note the same dependence of the complexity bound in Theorem 3.3 on $\epsilon$ as for the deterministic line-search algorithm, and the same dependence on $p$ as for the other problem classes discussed above.*

# 4 Probabilistic second-order models and cubic regularization methods

In this section we consider a randomized version of second-order methods, whose deterministic counterpart achieves optimal complexity rate [8, 5]. As in the line-search case, we show that in expectation, the same rate of convergence applies as in the deterministic (cubic regularization) case, augmented by a term that depends on the probability of having accurate models. Here we revert back to considering general objective functions that are not necessarily convex.

## 4.1 A cubic regularization algorithm with random models

Let us now consider a cubic regularization method where the following model

$$m_k(x^k + s) = f(x^k) + s^T g^k + \frac{1}{2} s^T b^k s + \frac{\sigma_k}{3} \|s\|^3, \tag{28}$$

is approximately minimized on each iteration $k$ with respect to $s$, for some vector $g^k$ and a matrix $b^k$ and some regularization parameter $\sigma^k > 0$. As before we assume that $g_k$ and $b_k$ are realizations of some random variables $G_k$ and $B_k$, which imply that the model is random and we assume that it is sufficiently accurate with probability at least $p$; the details of this assumption will be given after we state the algorithm.

The step $s^k$ is computed as in [7, 8] to approximately minimize the model (28), namely, it is required to satisfy

$$(s^k)^T g^k + (s^k)^T b^k s^k + \sigma_k \|s^k\|^3 = 0 \text{ and } (s^k)^T b^k s^k + \sigma_k \|s^k\|^3 \geq 0 \tag{29}$$

and

$$\|\nabla m_k(x^k + s^k)\| \leq \kappa_\theta \min\{1, \|s^k\|\} \|g^k\|, \tag{30}$$

where $\kappa_\theta \in (0,1)$ is a user-chosen constant.

Note that (29) is satisfied if $s^k$ is the global minimizer of the model $m_k$ over some subspace; in fact, it is sufficient for $s^k$ to be the global minimizer of $m_k$ along the line $\alpha s^k$ [8][2] Condition (30) is a relative termination condition for the model minimization (say over increasing subspaces) and it is clearly satisfied at stationary points of the model; ideally it will be satisfied sooner at least in the early iterations of the algorithm [8].

The probabilistic Adaptive Regularization with Cubics (ARC) framework is presented below.

## Algorithm 4.1 An ARC algorithm with random models

**Initialization**

Choose parameters $\gamma \in (0,1)$, $\theta \in (0,1)$, $\sigma_{\min} > 0$ and $\kappa_\theta \in (0,1)$. Pick initial $x^0$ and $\sigma_0 > \sigma_{\min}$. Repeat for $k = 0, 1, \ldots,$

---

[2]Note that a recently-proposed cubic regularization variant [2] can dispense with the approximate global minimization condition altogether while maintaining the optimal complexity bound of ARC. A probabilistic variant of [2] can be constructed similarly to probabilistic ARC, and our analysis here can be extended to provide same-order complexity bounds.

1. **Compute a model**
   Compute an approximate gradient $g^k$ and Hessian $b^k$ and form the model (28).

2. **Compute the trial step $s^k$**
   Compute the trial step $s^k$ to satisfy (29) and (30).

3. **Check sufficient decrease**
   Compute $f(x^k + s^k)$ and

$$\rho_k = \frac{f(x^k) - f(x^k + s^k)}{f(x^k) - m_k(x^k + s^k)}. \tag{31}$$

4. **Update the iterate**
   Set

$$x^{k+1} = \begin{cases} x^k + s^k & \text{if} \quad \rho_k \geq \theta \quad \quad [k \text{ successful}] \\ x^k & \text{otherwise} \quad \quad [k \text{ unsuccessful}] \end{cases} \tag{32}$$

5. **Update the regularization parameter $\sigma_k$**
   Set

$$\sigma_{k+1} = \begin{cases} \max\{\gamma\sigma_k, \sigma_{\min}\} & \text{if} \quad \rho_k \geq \theta \\ \frac{1}{\gamma}\sigma_k & \text{otherwise.} \end{cases} \tag{33}$$

**Remark 4.1** *Typically (see e.g. [7]) one would further refine (32) and (33) by distinguishing between successful and very successful iterations, when $\rho_k$ is not just positive but close to 1. It is beneficial in the deterministic setting to keep the regularization parameter unchanged on successful iterations when $\rho_k$ is greater than $\theta$ but is not close to 1 and only to decrease it when $\rho_k$ is substantially larger than $\theta$. For simplicity and uniformity of our general framework, we simplified the parameter update rule. However, the analysis presented here can be quite easily extended to the more general case by slightly extending the flexibility of the stochastic processes. In practice it is yet unclear if the same strategy will be beneficial, as "accidentally" bad models and the resulting unsuccessful steps may drive the parameter $\sigma_k$ to be larger than it should be, and hence a more aggressive decrease of $\sigma_k$ may be desired. This practical study is a subject of future research.*

**Remark 4.2** *We have stated Algorithm 4.1 so that it is as close as possible to known/deterministic ARC frameworks for ease of reading. We note however, that it is perfectly coherent with the generic algorithmic framework, Algorithm 2.1, if one sets $\alpha_k = 1/\sigma_k$ and $\rho_k \geq \theta$ as the sufficient decrease condition. We will exploit this connection in the analysis that follows.*

The requirement of sufficient model accuracy considered here is similar to the definition of probabilistically fully quadratic models introduced in [1], though note that we only require the second-order condition along the trial step $s^k$.

**Definition 4.1** *We say that a sequence of random models and corresponding directions $\{M_k\}$ is $(p)$-probabilistically "sufficiently accurate" for Algorithm 4.1 if there exist constants $\kappa_g$ and $\kappa_H$ such that for any corresponding random sequence $\{\mathcal{A}_k = 1/\Sigma_k, X^k\}$, the random indicator variables*

$$I_k = \mathbb{1}\{\|\nabla f(X^k) - G^k\| \leq \kappa_g\|S^k\|^2 \quad \text{and} \quad \|(H(X^k) - B^k)S^k\| \leq \kappa_H\|S^k\|^2\}$$

*satisfy the following submartingale-like condition*

$$P(I_k = 1|F_{k-1}^M) \geq p,$$

*where $F_{k-1}^M = \sigma(M_0, \ldots, M_{k-1})$ is the $\sigma$-algebra generated by $M_0, \ldots, M_{k-1}$.*

As before, for any realization of Algorithm 4.1, we refer to iterations $k$ when $I_k$ occurs as *true* iterations, and otherwise, as *false* iterations. Hence for all true iterations $k$,

$$\|\nabla f(x^k) - g^k\| \leq \kappa_g \|s^k\|^2 \quad \text{and} \quad \|(H(x^k) - b^k)s^k\| \leq \kappa_H \|s^k\|^2. \tag{34}$$

For the remainder of the analysis of Algorithm 4.1 we make the following assumption.

**Assumption 4.1** *The sequence of random models and corresponding directions $\{M_k, S_k\}$, generated in Algorithm 4.1, is $(p)$-probabilistically "sufficiently accurate" for the corresponding random sequence $\{\mathcal{A}_k = 1/\Sigma_k, X^k\}$, with $p > 1/2$.*

Regarding the possibly nonconvex objective $f$, in addition to Assumption 3.2, we also need the following assumption.

**Assumption 4.2** $f \in \mathcal{C}^2(\mathbb{R}^n)$ *and has globally Lipschitz continuous Hessian $H$, namely,*

$$\|H(x) - H(y)\| \leq L_H \|x - y\| \quad \text{for all } x, y \in \mathbb{R}^n \text{ and some } L_H > 0. \tag{35}$$

## 4.2 Global convergence rate analysis, nonconvex case

The next four lemmas give useful properties of Algorithm 4.1 that are needed later for our stochastic analysis.

**Lemma 4.1 (Lemma 3.3 in [7])** *Consider any realization of Algorithm 4.1. Then on each iteration $k$ we have*

$$f(x^k) - m_k(x^k + s^k) \geq \frac{1}{6}\sigma_k \|s^k\|^3. \tag{36}$$

*Thus on every successful iteration $k$, we have*

$$f(x^k) - f(x^{k+1}) \geq \frac{\theta}{6}\sigma_k \|s^k\|^3. \tag{37}$$

*Proof.* Clearly, (37) follows from (36) and the sufficient decrease condition (31)-(32). It remains to prove (36). Combining the first condition on step $s^k$ in (29), with the model expression (28) for $s = s^k$ we can write

$$f(x^k) - m_k(x^k + s^k) = \frac{1}{2}(s^k)^T B^k s^k + \frac{2}{3}\sigma_k \|s^k\|^3.$$

The second condition on $s^k$ in (29) implies $(s^k)^T B^k s^k \geq -\sigma_k \|s^k\|^3$, which, when used with the above equation, gives (36). $\qquad\square$

**Lemma 4.2** *Let Assumptions 3.2 and 4.2 hold. For any realization of Algorithm 4.1, if iteration $k$ is true (i.e., $I_k = 1$), and if*

$$\sigma_k \geq \sigma_c = \frac{2\kappa_g + \kappa_H + L + L_H}{1 - \frac{1}{3}\theta}, \tag{38}$$

*then iteration $k$ is also successful.*

*Proof.* Clearly, if $\rho_k - 1 \geq 0$, then $k$ is successful by definition. Let us consider the case when $\rho_k < 1$; then if $1 - \rho_k \leq 1 - \theta$, $k$ is successful. We have from (31), that

$$1 - \rho_k = \frac{f(x^k + s^k) - m_k(x^k + s^k)}{f(x^k) - m_k(x^k + s^k)}.$$

Taylor expansion and triangle inequalities give, for some $\xi^k \in [x^k, x^k + s^k]$,

$$
\begin{aligned}
&f(x^k + s^k) - m_k(x^k + s^k) \\
&= [\nabla f(x^k) - g^k]^T s^k + \tfrac{1}{2}(s^k)^T [H(\xi^k) - H(x^k)]s^k + \tfrac{1}{2}(s^k)^T [H(x^k) - b^k]s^k - \tfrac{1}{3}\sigma_k \|s^k\|^3 \\
&\leq \|\nabla f(x^k) - g^k\| \cdot \|s^k\| + \tfrac{1}{2}\|H(\xi^k) - H(x^k)\| \cdot \|s^k\|^2 + \tfrac{1}{2}\|(H(x^k) - b^k)s^k\| \cdot \|s^k\| - \tfrac{1}{3}\sigma_k \|s^k\|^3 \\
&\leq \left(\kappa_g + \tfrac{L_H}{2} + \tfrac{\kappa_H}{2} - \tfrac{1}{3}\sigma_k\right)\|s^k\|^3 = (6\kappa_g + 3L_H + 3\kappa_H - 2\sigma_k)\tfrac{1}{6}\|s^k\|^3,
\end{aligned}
$$

where the last inequality follows from the fact that the iteration is true and hence (34) holds, and from Assumption 4.2. This and (36) now give that $1 - \rho_k \leq 1 - \theta$ when $\sigma_k$ satisfies (38). $\square$

Note that for the above lemma to hold $\sigma_c$ does not have to depend on $L$. However, in what follows we will need another condition on $\sigma_c$, which will involve $L$; hence for simplicity of notation we introduced $\sigma_c$ above to satisfy all necessary bounds.

**Lemma 4.3** *Let Assumptions 3.2 and 4.2 hold. Consider any realization of Algorithm 4.1. On each true iteration $k$ we have*

$$\|s^k\| \geq \sqrt{\frac{1 - \kappa_\theta}{\sigma_k + \kappa_s}}\|\nabla f(x^k + s^k)\|, \tag{39}$$

*where $\kappa_s = 2\kappa_g + \kappa_H + L + L_H$.*

*Proof.* Triangle inequality, equality $\nabla m_k(x^k + s) = g^k + b^k s + \sigma_k \|s\| s$ and condition (30) on $s^k$ together give

$$
\begin{aligned}
\|\nabla f(x^k + s^k)\| &\leq \|\nabla f(x^k + s^k) - \nabla m_k(x^k + s^k)\| + \|\nabla m_k(x^k + s^k)\| \\
&\leq \|\nabla f(x^k + s^k) - g^k - b^k s^k\| + \sigma_k \|s^k\|^2 + \kappa_\theta \min\{1, \|s^k\|\}\|g^k\|.
\end{aligned} \tag{40}
$$

Recalling Taylor expansion of $\nabla f(x^k)$

$$\nabla f(x^k + s^k) = \nabla f(x^k) + \int_0^1 H(x^k + ts^k)s^k dt,$$

and applying triangle inequality, again, we have

$$
\begin{aligned}
\|\nabla f(x^k + s^k) - g^k - b^k s^k\| &\leq \|\nabla f(x^k) - g^k\| + \\
&\quad \left\| \int_0^1 [H(x^k + ts^k) - H(x^k)]s^k dt \right\| + \|H(x^k)s^k - b^k s^k\| \\
&\leq \left\{ \kappa_g + \tfrac{1}{2}L_H + \kappa_H \right\} \|s^k\|^2,
\end{aligned}
$$

where to get the second inequality, we also used (34) and Assumption 4.2.

We can bound $\|g^k\|$ as follows

$$
\|g^k\| \leq \|g^k - \nabla f(x^k)\| + \|\nabla f(x^k) - \nabla f(x^k + s^k)\| + \|\nabla f(x^k + s^k)\| \leq \kappa_g \|s^k\|^2 + L\|s^k\| + \|\nabla f(x^k + s^k)\|.
$$

Thus finally, we can bound all the terms on the right hand side of (40) in terms of $\|s^k\|^2$ and using the fact that $\kappa_\theta \in (0, 1)$ we can write

$$
(1 - \kappa_\theta)\|\nabla f(x^k + s^k)\| \leq (2\kappa_g + \kappa_H + L + L_H + \sigma_k)\|s^k\|^2,
$$

which is equivalent to (39). $\qquad \square$

**Lemma 4.4** *Let Assumptions 3.2 and 4.2 hold. Consider any realization of Algorithm 4.1. On each true and successful iteration $k$, we have*

$$
f(x^k) - f(x^{k+1}) \geq \frac{\kappa_f}{(\max\{\sigma_k, \sigma_c\})^{3/2}} \|\nabla f(x^{k+1})\|^{3/2}, \tag{41}
$$

*where $\kappa_f := \frac{\theta}{12\sqrt{2}}(1 - \kappa_\theta)^{3/2}\sigma_{\min}$ and $\sigma_c$ is defined in (38).*

*Proof.* Combining Lemma 4.3, inequality (37) from Lemma 4.1 and the definition of successful iteration in Algorithm 4.1 we have, for all true and successful iterations $k$,

$$
f(x^k) - f(x^{k+1}) \geq \frac{\theta}{6}(1 - \kappa_\theta)^{3/2} \frac{\sigma_k}{(\sigma_k + \kappa_s)^{3/2}} \|\nabla f(x^{k+1})\|^{3/2}. \tag{42}
$$

Using that $\sigma_k \geq \sigma_{\min}$ and that $\kappa_s \leq \sigma_c$, (42) implies (41). $\qquad \square$

**The stochastic processes and global convergence rate analysis**  We are now ready to cast Algorithm 4.1 and its behavior into the generic stochastic analysis framework of Section 2.

For each realization of Algorithm 4.1, we define

$$
\alpha_k = \frac{1}{\sigma_k} \quad \text{and} \quad f_k = f(x^0) - f(x^k),
$$

and consider the corresponding stochastic process $\{\mathcal{A}_k = 1/\Sigma_k, F_k = f(X^0) - f(X^k)\}$. Let $F_\epsilon = f(x^0) - f^*$ denote the upper bound on the progress measure $F_k$.

As in the case of the line-search algorithm applied to nonconvex objectives, we would like to bound the expected number of iterations that Algorithm 4.1 takes until $\|\nabla f(X^k)\| \leq \epsilon$ occurs. Here, however, for technical reasons made clear below, we count the number of iterations until a successful iteration results in $x^{k+1}$ such that $\|\nabla f(X^{k+1})\| \leq \epsilon$. Let $N_\epsilon$ denote the (random) index of such an iteration. (Clearly, $N_\epsilon$ thus defined is simply one less than the number of iterations that occur until $\|\nabla f(X^{k+1})\| \leq \epsilon$.)

Regarding Assumption 2.1, Lemmas 4.2 and 4.4 provide that the following must hold for any realization of Algorithm 4.1.

- If $k$ is a true and successful iteration, then

$$f_{k+1} \geq f_k + \frac{\kappa_f}{(\max\{\sigma_k, \sigma_c\})^{3/2}} \|\nabla f(x^{k+1})\|^{3/2}$$

  and

$$\alpha_{k+1} = \gamma^{-1}\alpha_k.$$

- If $\alpha_k \leq C = \frac{1}{\sigma_c}$, where $\sigma_c$ is defined in (38), and iteration $k$ is true, then it is also successful.

Hence, once again, Assumption 2.1 holds and the process $\{\mathcal{A}_k, F_k\}$ behaves exactly as our generic process (2)-(3) in Section 2.4, with $C = \frac{1}{\sigma_c} = \frac{1 - \frac{1}{3}\theta}{2\kappa_g + \kappa_H + L + L_H}$, and the specific choice

$$h(\mathcal{A}_k) = \kappa_f(\min\{\mathcal{A}_k, C\})^{3/2}\epsilon^{3/2}.$$

for all $k < N_\epsilon$.

Finally, the complexity result again follows from Theorem 2.1 and the expressions for $C$, $h(C)$ and $F_\epsilon$.

**Theorem 4.1** *Let Assumptions 3.2, 4.1 and 4.2 hold. Then the expected number of iterations that Algorithm 4.1 takes until $\|\nabla f(X^{k+1})\| \leq \epsilon$ occurs is bounded by*

$$\mathbb{E}(N_\epsilon) \leq \frac{2p}{(2p-1)^2}\left(\frac{M}{\epsilon^{3/2}} + \log\left(\frac{2\kappa_g + \kappa_H + L + L_H}{\sigma_0(1 - 1/3\theta)}\right)\right),$$

*where $M = \frac{(f(x^0) - f^*)(2\kappa_g + \kappa_H + L + L_H)^{3/2}}{\kappa_f(1 - 1/3\theta)^{3/2}}$ is a constant independent of $p$ and $\epsilon$.*

**Remark 4.3** *We note that the dependency on $\epsilon$ in the above bound on the expected number of iterations is of the order $\epsilon^{-3/2}$, which is of the same order as for the deterministic ARC algorithm and is the optimal rate for nonconvex optimization using second order models [5]. The dependence on $p$ is, again, the same as in the case of line-search and it is intuitive.*

**Remark 4.4** *Theorem ?? stating that $\liminf k \to \infty \|\nabla f(X^k)\| = 0$ almost surely, holds for Algorithm 4.1 since a similar proof applies.*

## 5 Random models

In this section we will discuss and motivate the definition of probabilistically "sufficiently accurate" models. In particular, Definition 3.1 is a modification of the definition of probabilistically fully-linear models, which is used in [1]. Similarly, Definition 4.1 is similar to that of probabilistically fully-quadratic models in [1]. These definitions serve to provide properties of the model (with some probability) which are sufficient for first-order (in the case of Definition 3.1) and second-order (in the case of Definition 4.1) convergence rates.

We will now describe several setting where the models are random and satisfy our definitions.

## 5.1 Stochastic gradients and batch sampling

In [4] an adaptive sample size strategy was proposed in the setting where $\nabla f(x) = \sum_{i=1}^{N} \nabla f_i(x)$, for large values of $N$. In this case computing $\nabla f(x)$ accurately can be prohibitive, hence, instead an estimate $\nabla f_S(x) = \sum_{i \in S} \nabla f_i(x)$ is often computed in hopes that it provides a good estimate of the gradient and a descent direction. It is observed in [4] that if sample sets $S_k$ on each iteration ensure that

$$\|\nabla f_{S_k}(x^k) - \nabla f(x^k)\| \leq \mu \|\nabla f_{S_k}(x^k)\| \tag{43}$$

for some $\mu \in (0,1)$, then using a fixed step size

$$\alpha_k \equiv \alpha \leq \frac{1-\mu}{L} \tag{44}$$

the step $s_k = -\alpha \nabla f_{S_k}(x)$ is always a descent step and the line search algorithm converges with the rate $O(\log(1/\epsilon))$ if $f$ is strongly convex. Clearly, condition (43) implies that the model $M_k(x) = f(x^k) + \nabla f_{S_k}(x^k)^\top (x - x^k)$ is sufficiently accurate according to Definition 3.1 for the given fixed step size $\alpha$. Hence Assumption 3.1 on the models can be viewed as a relaxed version of those in [4], since we allow the condition (43) to fail, as long as it fails with probability less than $1/2$, conditioned on the past. Moreover, we analyze the practical version of line search algorithm, with a variable step size, which does not have to remain smaller than $\frac{1-\mu}{L}$ and we provide convergence rates in convex, strongly convex and nonconvex setting.

Convergence in expectation of a stochastic algorithm is further shown in [4]. In particular, under the assumption that the variance of $\|\nabla f_i(x)\|$ is bounded for all $i$ and that $\mathbb{E}_S[\nabla f_S(x^k)] = \nabla f(x^k)$, it is shown that, for $X^k$ computed after $k$ steps of stochastic gradient descent with a fixed step size, $\mathbb{E}[f(X^k)]$ converges linearly of $f^*$, when $f(x)$ is strongly convex and if $|S_k|$ - the size of the sample set $S_k$ - grows exponentially with $k$.

Here, again, our results can be viewed as a generalization of the results in [4]. Indeed, let us assume that $\mathbb{E}_S[\nabla f_S(x^k)] = \nabla f(x^k)$ for each $x^k$ and let $t_k = |S_k|$ - the size of the sample set $S_k$. Since variance of $\|\nabla f_i(x)\|$ is bounded for all $i$, we have that $\mathbb{E}_{S_k}[\|\nabla f_{S_k}(x^k) - \nabla f(x^k)\|] \leq \frac{w}{t_k}$, for some fixed $w$, where the expectation is taken over all random sample sets $S_k$ of size $t_k$. In other words, the variance of one sample of the stochastic gradient $\|\nabla f_i(x)\|$ is bounded and hence the variance of $\nabla f_{S_k}(x^k)$ decays as the size of $S_k$ increases.

By Chebychev inequality

$$Pr\{\|\nabla f_{S_k}(x^k) - \nabla f(x^k)\| > \min\{1/2, \alpha_k\} \|\nabla f(x^k)\|\} \leq \frac{w}{\min\{1/2, \alpha_k\}^2 \|\nabla f(x^k)^2\| |S_k|}.$$

If $\|\nabla f_{S_k}(x^k) - \nabla f(x^k)\| \leq \min\{1/2, \alpha_k\} \|\nabla f(x^k)\|$, for a particular $x_k$ and a sample set $S_k$, then by applying triangle inequality we have

$$\|\nabla f_{S_k}(x^k) - \nabla f(x^k)\| \leq \frac{\alpha_k \|\nabla f_{S_k}(x^k)\|}{2}.$$

Hence the probability of the event $\|\nabla f_S(x^k) - \nabla f(x^k)\| \leq \frac{\alpha_k \|\nabla f_{S_k}(x^k)\|}{2}$ is at least

$$1 - \frac{w}{\min\{1/2, \alpha_k\}^2 \|\nabla f(x^k)\|^2 |S_k|} \geq 1 - \frac{w}{\min\{1/2, \alpha_k\}^2 (1 + \alpha_k)^2 \|\nabla f_{S_k}(x^k)\|^2 |S_k|},$$

hence as long as $|S_k|$ is chosen sufficiently large, then this probability is greater than $1/2$ and $\nabla f_S(x^k)$ provides us with a probabilistically sufficiently accurate model according to Definition

27

3.1. Hence the theory described in this paper applies to the case of line search based on stochastic gradient. Note that, on top of the results in [4], we not only analyze line search in nonconvex and convex setting, but also show the bound on the expected number of iterations until the desired accuracy is reached, rather than the expected accuracy after a given number of iterations. As we have shown earlier, this implies lim inf-type convergence with probability one. Moreover, as mentioned on page 13, it is not difficult to extend our analysis to show that the number of iterations until the desired accuracy is reached has exponentially decaying tails.

Analyzing complexity of methods in this setting in terms of the total number of gradient samples is a subject of some current research [18]. We leave the exact comparison that can be obtained from our results and those existing in current literature as future research, as this requires defining a sample size selection strategy and possible improvement of our results. Similarly, we leave for future research the derivations of the models in this setting that satisfy Definition 4.1 for the use within the ARC algorithm.

## 5.2   Models based on random sampling of function values

The motivation behind the notions of probabilistically fully-linear and fully-quadratic models introduced in [1] is based on derivative-free models, which are models based on function values, rather than gradient estimates. We will now show how such models fit into our framework.

Let us first recall the definition of probabilistically fully-linear and quadratic models and pose it in the terms closest to the ones used in this paper

**Definition 5.1**    *1. We say that a sequence of random models $\{M_k\}$ is $(p)$-probabilistically fully-linear if there exists constant $\kappa_g$ such that for any corresponding random sequence $\Delta_k$, $X^k$, the random indicator variables*

$$I_k^l \ = \ \mathbb{1}\left\{\|\nabla f(X^k) - G^k\| \leq \kappa_g \Delta_k\right\}$$

*satisfy the following submartingale-like condition*

$$P(I_k^l = 1 | F_{k-1}^M) \ \geq \ p,$$

*where $F_{k-1}^M = \sigma(M_0, \ldots, M_{k-1})$ is the $\sigma$-algebra generated by $M_0, \ldots, M_{k-1}$.*

*2. We call sequence $\{M_k\}$ is $(p)$-probabilistically fully-quadratic if there exist constants $\kappa_g$ and $\kappa_H$ such that for any corresponding random sequence $\Delta_k$, $X^k$, the random indicator variables*

$$I_k^q \ = \ \mathbb{1}\{\|\nabla f(X^k) - G^k\| \leq \kappa_g \Delta_k^2 \quad \text{and} \quad \|H(X^k) - B^k\| \leq \kappa_H \Delta_k\}$$

*satisfy the following submartingale-like condition*

$$P(I_k^q = 1 | F_{k-1}^M) \ \geq \ p,$$

*where $F_{k-1}^M = \sigma(M_0, \ldots, M_{k-1})$ is the $\sigma$-algebra generated by $M_0, \ldots, M_{k-1}$.*

The key difference between the conditions in Definition 5.1 and those in Definitions 3.1 and 4.1 is the right hand side of the error bounds - in the case of fully-linear and fully-quadratic models $\Delta_k$ is a random variable that does not depend on $M_k$, but in the case of this paper, $\Delta_k$

is replaced by $\mathcal{A}_k \|G_k\|$ in the case of Definition 3.1 and by $\|S_k\|$ in the case of Definition 4.1. In other words, the accuracy of the model has to be proportional to the step size which this model produces. Since in [1] trust region methods are analyzed instead of line search and ARC, Definition 5.1 is sufficient.

Models in [1] are constructed by sampling function values in a ball of a given radius around the current iterate $x^k$ and in all cases construction of the $k$-th model $M_k$ relies on the knowledge of the sampling radius. We will now show that, given a mechanism of constructing probabilistically fully-linear and fully-quadratic models for any sequence of radii (as described in [1]), we can modify our line search algorithm and ARC algorithm, respectively, and extend the convergence rate analysis to utilize these models.

**Line-search with probabilistically fully-linear models**   Let us consider Algorithm 3.1 and corresponding random sequence of iterates $X^k$ and step sizes $\mathcal{A}_k$. If a given model $M_k$ is fully-linear in $B(X^k, \mathcal{A}_k \Xi_k)$ and $\|G_k\| \geq \kappa_\Delta \Xi_k$, for some positive constant $\kappa_\Delta$, then model $M_k$ is sufficiently accurate, according to Definitions 3.1.

To achieve this, for instance, in nonconvex case, for all $\|\nabla f(X^k)\| \geq \epsilon$ consider $\Xi_k \leq \frac{\epsilon}{2\kappa_g \max\{\mathcal{A}_k, 1\}}$, where $\kappa_g$ is the constant in the definition of fully-linear models. Then any fully-linear model $M_k(x)$ is also sufficiently accurate, simply because $\|\nabla f(X^k) - G^k\| \leq \kappa_g \mathcal{A}_k \Xi_k \leq \min\{\mathcal{A}_k, 1\} \frac{\epsilon}{2}$ implies $\|G_k\| \geq \frac{\epsilon}{2} \geq \Xi_k$. Similar bounds can be derived for the convex and strongly convex cases.

Consider the following example of a method that produces probabilistically sufficiently accurate models, based on the arguments above. Suppose we are estimating gradients of $f(x)$ by a finite difference scheme using step size $\mathcal{A}_k \Xi_k$, with $\Xi_k$, sufficiently small, and suppose we compute the function values using parallel computations. If some of the computations fail to complete (due to an overloaded processor, say) with some probability and the total probability of having a computational failure in any of the processors at each iteration is less than $1/2$, conditioned on the past, then we obtain probabilistically sufficiently accurate models. Note that, we do not assume the nature of the computational error, when such error occurs, hence allowing for the gradient estimate to be, occasionally, completely inaccurate.

Another example can be derived from [1], where it is shown that sparse gradient and Hessian estimates can be obtained by randomly sampling fewer function values than is needed to construct gradient and/or Hessian by finite differences. Using this sampling strategy, probabilistically fully-linear and fully-quadratic models can be generated at reduced computation cost. Here again, choosing sampling radius to equal $\Delta_k = \mathcal{A}_k \Xi_k$, with sufficiently small $\Xi_k$ will guarantee that the models are also probabilistically sufficiently accurate.

We now address a more practical approach, when estimates $\Xi_k$ are not chosen to be small enough a priory, but are dynamically decreased, as another parameter in the algorithm. We will outline how our theory can be extended in this case. Consider the following modification of Algorithm 3.1.

**Algorithm 5.1 Line-search with probabilistically fully-linear models**

**Initialization**

   *Chose constants $\theta \in (0,1)$, $\gamma \in (0,1)$, $\alpha_{\max} > 0$ and $\kappa_\Delta > 1$. Pick initial $x^0$ and $\alpha_0 < \alpha_{\max}$, $\xi_0$. Repeat for $k = 0, 1, \ldots$*

1. **Compute a model**

   *Compute a model $m_k$, which is probabilistically fully-linear in $B(x^k, \alpha_k \xi_k)$ and use it to generate a direction $g^k$.*

2. **Check model accuracy**

   *If $\|g^k\| \geq \kappa_\Delta \xi_k$, then set the step $s^k = -\alpha_k g^k$ and continue to Step 3.*
   *Otherwise, $x^{k+1} = x^k$, $\alpha_{k+1} = \alpha_k$, $\xi_{k+1} = \xi_k/\kappa_\Delta$, return to Step 1.*

3. **Check sufficient decrease**

   *Check if*
   $$f(x^k - \alpha_k g^k) \leq f(x^k) - \alpha_k \theta \|g^k\|^2. \tag{45}$$

4. **Successful step**

   *If (45) holds, then $x^{k+1} := x^k - \alpha_k g^k$ and $\alpha_{k+1} = \min\{\alpha_k/\gamma, \alpha_{\max}\}$.*

5. **Unsuccessful step**

   *Otherwise, $x^{k+1} = x^k$.*
   $\alpha_{k+1} = \gamma \alpha_k$.

In the above algorithm, at each iteration we maintain $\xi_k$, which is expected to be an underestimate of the norm of the descent direction, up to a constant, $\kappa_\Delta$. The algorithm then uses $\delta_k = \alpha_k \xi_k$ as the radius for constructing fully-linear models. After the model is produced, condition $\|g_k\| \geq \kappa_\Delta \xi_k$ is checked. If this condition holds, the algorithm proceeds exactly as the original version, but if this condition fails, then $\xi_k$ is reduced by a constant ($\kappa_\Delta$ is a practical choice, but any other constant can be used) and the iteration is declared to be unsuccessful (hence $x^{k+1} = x^k$), and the step size $\alpha_k$ remains the same.

Let us consider different possible outcomes for each iteration $k$ for which $\|\nabla f(x^k)\| \geq \epsilon$. From our analysis above, we know that if $\xi_k \leq \frac{\epsilon}{2\kappa_g \alpha_{\max}}$ and the model $m_k$ is fully linear, then $\|g_k\| \geq \xi_k$, hence the model is also sufficiently accurate and the iteration of Algorithm 5.1 proceeds as in Algorithm 3.1. Since $\xi_k$ is never increased, then, once it is small enough, the analysis of Algorithm 5.1 can be reduced to that of Algorithm 3.1. Then what remains is to estimate the number of iterations that Algorithm 5.1 takes until $\xi_k \leq \frac{\epsilon}{2\kappa_g \alpha_{\max}}$ or $\|\nabla f(x^k)\| \leq \epsilon$ occurs.

While $\xi_k$ is not sufficiently small, we can have the following outcomes: 1) $\|g_k\| < \kappa_\Delta \xi_k$, in which case $\xi_k$ is reduced, 2) the model is not fully linear and $\|g_k\| \geq \kappa_\Delta \xi_k$, hence the model may not be sufficiently accurate, but $\xi_k$ is not reduced and 3) the model is fully-linear and $\|g_k\| \geq \kappa_\Delta \xi_k$, hence the model is also sufficiently accurate. Hence with probability at least $p$, $\xi_k$ is reduced or the model is sufficiently accurate. It is possible to extend the definition of our stochastic processes and their analysis to compute the upper bounds on the expected number of iterations Algorithm 5.1 takes until $\|\nabla f(x^k)\| \leq \epsilon$ occurs. This bound will be increased by adding a constant times the number of iterations it takes to achieve $\xi_k \leq \frac{\epsilon}{2\kappa_g \alpha_{\max}}$, which is $O(\log(1/\epsilon))$. Again, similar analysis can be carried out for the cases of convex and strongly convex functions.

**ARC with probabilistically fully-quadratic models** Let us consider Algorithm 4.1. In this case, in the same vein with line-search, we consider setting $\Delta_k$ in the Definition 5.1 of probabilistically fully-quadratic models, to a sufficiently small value or adjusting it in the run of the algorithm so as to ensure that when the model is fully-quadratic, it is also sufficiently

accurate (at least asymptotically). We will make these two approaches to the choice of $\Delta_k$ more precise in what follows. To this end, we need a new variant of Lemma 4.3 for the case of probabilistically fully-quadratic models.

**Lemma 5.1** *Let Assumptions 3.2 and 4.2 hold. Consider any realization of Algorithm 4.1 where we generate models that are p-probabilistically fully-quadratic according to Definition 5.1. Then on each iteration $k$ in which $I_k^q = 1$, we have*

$$(1 - \kappa_\theta)\|\nabla f(x^k + s^k)\| \leq (2\kappa_g + \kappa_H)\delta_k \max\{\delta_k, 1\} + (L + L_H + \sigma_k)\|s^k\|^2. \tag{46}$$

*In particular, if $\epsilon \in (0, 1]$, $\max\{L, L_H\} \geq 1$, and*

$$\delta_k \leq \frac{(1 - \kappa_\theta)\epsilon}{\max\left\{2(2\kappa_g + \kappa_H), L + L_H + \sigma_k\right\}}, \tag{47}$$

*then on each iteration $k$ with $\|\nabla f(x^k + s^k)\| \geq \epsilon$ and in which $I_k^q = 1$, we have $\|s^k\| \geq \delta_k$.*

*Assume now that $\delta_k = \frac{\xi_k}{\sigma_k}$. Then if $\epsilon \in (0, 1]$, $\max\{L, L_H\} \geq 1$, and*

$$\xi_k \leq \frac{(1 - \kappa_\theta)\sigma_{\min}\epsilon}{\max\left\{2(2\kappa_g + \kappa_H), L + L_H + \sigma_{\min}\right\}} := \xi_\epsilon, \tag{48}$$

*then on each iteration $k$ with $\|\nabla f(x^k + s^k)\| \geq \epsilon$ and in which $I_k^q = 1$, we have $\|s^k\| \geq \delta_k$.*

*Proof.* It follows from Definition 5.1 that on each realization of Algorithm 4.1, we have

$$\|\nabla f(x^k) - g^k\| \leq \kappa_g \delta_k^2 \quad \text{and} \quad \|H(x^k) - b^k\| \leq \kappa_H \delta_k \tag{49}$$

The proof of (46) now follows identically to the proof of Lemma 4.3 if one uses (49) instead of (34).

The choice of $\delta_k$ in (47) implies $\delta_k \leq 1$ and so $\|s^k\| \geq \delta_k$ trivially holds when $\|s^k\| \geq 1$. When $\|s^k\| < 1$, $\|\nabla f(x^k + s^k)\| \geq \epsilon$, and $\delta_k \leq 1$, (46) implies

$$(1 - \kappa_\theta)\epsilon - (2\kappa_g + \kappa_H)\delta_k \leq (L + L_H + \sigma_k)\|s^k\|.$$

Now the condition (47) on $\delta_k$ implies $(L + L_H + \sigma_k)\|s^k\| \geq (1 - \kappa_\theta)\epsilon/[2(2\kappa_g + \kappa_H)]$. Applying again the upper bound on $\delta_k$ provides $\|s^k\| \geq \delta_k$.

Finally, if $\delta_k = \frac{\xi_k}{\sigma_k}$, and using $\sigma^k \geq \sigma_{\min}$ for all $k$ due to the algorithm construction, (48) implies (47). $\square$

The second part of Lemma 5.1 provides that if p-probabilistically fully-quadratic models are generated with $\delta_k$ chosen sufficiently small so that (47) holds, then the models are also p-probabilistically sufficiently accurate. Thus Algorithm 4.1 can be run with models sampled in this way and the analysis carries through as before. For example, as in the case of linesearch, $g^k$ and $b^k$ could be generated by (sufficiently accurate) finite-difference schemes using function values, where computations are done in parallel and where the total probability of computational failure in any of the processors at each iteration is less than $1/2$.

Note however, that the bound that dictates the choice of a suitably small $\delta_k$ depends on problem constants that may not be known a priori. Thus it would be better – and computationally more efficient – to adjust $\delta_k$ during the run of Algorithm 4.1. A modification of Algorithm 4.1 that allows this is given next, and can be viewed as the analogue for ARC of the line-search Algorithm 5.1.

**Algorithm 5.2 ARC with probabilistically fully-quadratic models**

**Initialization**

> *Choose parameters $\sigma_{\min} > 0$, $\gamma \in (0,1)$, $\theta \in (0,1)$, $0 < \kappa_\theta < 1$ and $\kappa_\Delta > 1$. Pick a starting point $x^0$, a starting value $\sigma_0 > \sigma_{\min}$ and $\xi_0 > 0$. Repeat for $k = 0, 1, \ldots,$*

**1. Compute a model**

> *Compute a model which is probabilistically fully-quadratic in $B\left(x^k, \frac{\xi_k}{\sigma_k}\right)$, and hence generate approximate gradient $g^k$ and Hessian $b^k$.*

**2. Compute the trial step $s^k$**

> *Compute the trial step $s^k$ to satisfy (29) and (30).*

**3. Check model accuracy**

> *If $\|s^k\| \geq \kappa_\Delta \delta_k := \xi_k / \sigma_k$, then go to Step 4.*
> *Otherwise, set $x^{k+1} = x^k$, $\sigma_{k+1} = \sigma_k$, $\xi_{k+1} = \xi_k / \kappa_\Delta$, and return to Step 1.*

**4. Check sufficient decrease**

> *Compute $f(x^k + s^k)$ and*

$$\rho_k = \frac{f(x^k) - f(x^k + s^k)}{f(x^k) - m_k(x^k + s^k)}.$$

**5. Update the iterate**

> *Set*

$$x^{k+1} = \begin{cases} x^k + s^k & \text{if} \quad \rho_k \geq \theta \quad\quad [k \text{ successful}] \\ x^k & \text{otherwise} \quad\quad\quad\ [k \text{ unsuccessful}] \end{cases}$$

**6. Update the regularization parameter $\sigma_k$**

> *Set*

$$\sigma_{k+1} = \begin{cases} \max\{\gamma \sigma_k, \sigma_{\min}\} & \text{if} \quad \rho_k \geq \theta \\ \frac{1}{\gamma} \sigma_k & \text{otherwise.} \end{cases}$$

Algorithm 5.2 updates $\xi_k$ in order to obtain an underestimate $\delta_k := \xi_k / \sigma_k$ on the length of the step $s^k$. It constructs probabilistically fully-quadratic models in $B(x^k, \xi_k / \sigma_k)$ and checks whether $\|s^k\| \geq \kappa_\Delta \delta_k$. If that is the case, then the iteration of the above algorithm proceeds as (Algorithm 4.1) before; note that then, if the model is fully quadratic then it is also sufficiently accurate. Otherwise, if the step is too short, then $\xi_k$ is decreased by $\kappa_\Delta$, $x^k$ and $\sigma^k$ remain unchanged and a new model is generated (within the smaller ball).

Let us consider the behavior of Algorithm 5.2 while $\|\nabla f(x^k + s^k)\| \geq \epsilon$. It follows from the last part of Lemma 5.1 that since $\xi_\epsilon$ is independent of $k$ and $\xi_k$ is never increased in the algorithm, if $\kappa_\Delta \xi_j \leq \xi_\epsilon$ for some $j$, then $\xi_k$ will remain below this threshold for all subsequent iterations $k \geq j$; from this $j$ onwards, whenever the model is fully quadratic, then $\|s^k\| \geq \kappa_\Delta \delta_k$ and the model is also sufficiently accurate. Thus from iteration $j$ onwards, Algorithm 5.2 reduces to Algorithm 4.1 and the complexity analysis is the same as before. It remains to estimate the size of $j$, namely, the number of iterations Algorithm 5.2 takes until $\xi_k \leq \xi_\epsilon$ or $\|\nabla f(x^k + s^k)\| < \epsilon$.

Similarly to the linesearch analysis of possible outcomes above, we can argue that while $\xi_k$ is not sufficiently small, at least with probability $p$, $\xi_k$ is reduced or the model is sufficiently accurate. Thus, extending our earlier ARC analysis (and definitions of stochastic processes, etc)

to account for the $\xi_k$ updates as well, we would find that the complexity bound for Algorithm 5.2 is essentially that of Algorithm 4.1 plus a $\mathcal{O}(\log(1/\epsilon))$ term (coming from $\log(\xi_0/\xi_\epsilon)\log\kappa_\Delta$) that accounts for the number of iterations to drive $\xi_k$ below $\xi_\epsilon$.

# 6  Conclusions

We have proposed a general algorithmic framework with random models and a methodology for analyzing its complexity that relies on bounding the hitting time of a nondecreasing stochastic process that measures progress towards optimality. Our framework accounts for linesearch and cubic regularization methods, for example, and we particularize our results to obtain precise complexity bounds in the case of nonconvex and convex functions. Despite allowing our models to be arbitrarily inaccurate sometimes, the bounds we obtained match their deterministic counterparts in the order of the accuracy $\epsilon$. The effect of model inaccuracy is reflected by the constant multiple of the bound, which is a function of the probability that the model is sufficiently accurate. We have also briefly discussed ways to obtain probabilistically sufficiently accurate models as required by our framework.

The results in the paper assume that the objective $f$ is deterministic. Obtaining global rates of convergence results for similar algorithmic frameworks when $f$ is stochastic is a topic of future research. Also, further exploring ways to efficiently generate probabilistically sufficiently accurate models may increase the applicability of our results to a diverse set of problems.

# References

[1] A. BANDEIRA, K. SCHEINBERG, AND L. VICENTE, *Convergence of trust-region methods based on probabilistic models*, SIAM Journal on Optimization, 24 (2014), pp. 1238–1264.

[2] E. G. BIRGIN, J. L. GARDENGHI, S. A. S. J. M. MARTINEZ, AND P. L. TOINT, *Worst-case evaluation complexity for unconstrained nonlinear optimization using high-order regularized models*, Tech. Rep. naXys-05-2015, Department of Mathematics, University of Namur, 2015.

[3] R. BYRD, J. NOCEDAL, AND F. OZTOPRAK, *An inexact successive quadratic approximation method for convex l-1 regularized optimization*, tech. rep., 2013.

[4] R. H. BYRD, G. M. CHIN, J. NOCEDAL, AND Y. WU, *Sample size selection in optimization methods for machine learning*, Math. Program., 134 (2012), pp. 127–155.

[5] C. CARTIS, N. GOULD, AND P. L. TOINT, *Optimal Newton-type methods for nonconvex smooth optimization problems*, Tech. Rep. Optimization Online, 2011.

[6] ——, *On the oracle complexity of first-order and derivative-free algorithms for smooth nonconvex minimization*, SIAM Journal on Optimization, 22 (2012), pp. 66–86.

[7] C. CARTIS, N. I. M. GOULD, AND P. L. TOINT, *Adaptive cubic regularisation methods for unconstrained optimization. Part I: motivation, convergence and numerical results*, Math. Program., 127 (2011), pp. 245–295.

[8] ———, *Adaptive cubic regularisation methods for unconstrained optimization. Part II: worst-case function- and derivative-evaluation complexity*, Math. Program., 130 (2011), pp. 295–319.

[9] R. CHEN, *Stochastic Derivative-Free Optimization of Noisy Functions*, PhD thesis, Department of Industrial and Systems Engineering, Lehigh University, Bethlehem, USA, 2015.

[10] R. CHEN, M. MENICKELLY, AND K. SCHEINBERG, *Stochastic optimization using a trust-region method and random models*, tech. rep., ISE Dept., Lehigh University.

[11] O. DEVOLDER, F. GLINEUR, AND Y. NESTEROV, *First-order methods of smooth convex optimization with inexact oracle*, Math. Program., 146 (2014), pp. 37–75.

[12] S. GHADIMI AND G. LAN, *Stochastic first- and zeroth-order methods for nonconvex stochastic programming*, SIAM Journal on Optimization, 23 (2013), pp. 2341–2368.

[13] S. GRATTON, C. W. ROYER, L. N. VICENTE, AND Z. ZHANG, *Direct search based on probabilistic descent*, Tech. Rep. 14-11, Dept. Mathematics, Univ. Coimbra, 2014.

[14] J. D. LEE, Y. SUN, AND M. A. SAUNDERS, *Proximal newton-type methods for convex optimization*, in NIPS, 2012.

[15] Y. NESTEROV, *Introductory Lectures on Convex Optimization*, Kluwer Academic Publishers, Dordrecht, The Netherlands, 2004.

[16] Y. NESTEROV, *Random gradient-free minimization of convex functions*, Tech. Rep. 2011/1, CORE, 2011.

[17] Y. NESTEROV AND B. T. POLYAK, *Cubic regularization of newton method and its global performance*, Math. Program., 108 (2006), pp. 177–205.

[18] R. PASUPATHY, P. W. GLYNN, S. GHOSH, AND F. HAHEMI, *How much to sample in simulation-based stochastic recursions?*, (2014). Under Review.

[19] H. ROBBINS AND S. MONRO, *A stochastic approximation method*, Annals of Mathematical Statistics, 22 (1951), pp. 400–407.

[20] M. W. SCHMIDT, N. L. ROUX, AND F. BACH, *Convergence rates of inexact proximal-gradient methods for convex optimization*, in NIPS, 2011, pp. 1458–1466.

[21] ———, *Minimizing finite sums with the stochastic average gradient*, CoRR, abs/1309.2388 (2013).

[22] A. SHIRYAEV, *Probability*, Graduate Texts on Mathematics, Springer-Verlag, New York, 1995.

[23] J. SPALL, *Multivariate stochastic approximation using a simultaneous perturbation gradient approximation*, IEEE Transactions on Automatic Control, 37 (1992), pp. 332–341.