# Sample Complexity Analysis for Adaptive Optimization Algorithms with Stochastic Oracles

Billy Jin[*]        Katya Scheinberg[*]        Miaolan Xie[*]

October 2, 2023

### Abstract

Several classical adaptive optimization algorithms, such as line search and trust-region methods, have been recently extended to stochastic settings where function values, gradients, and Hessians in some cases, are estimated via stochastic oracles. Unlike the majority of stochastic methods, these methods do not use a pre-specified sequence of step size parameters, but adapt the step size parameter according to the estimated progress of the algorithm and use it to dictate the accuracy required from the stochastic oracles. The requirements on the stochastic oracles are, thus, also adaptive and the oracle costs can vary from iteration to iteration. The step size parameters in these methods can increase and decrease based on the perceived progress, but unlike the deterministic case they are not bounded away from zero due to possible oracle failures, and bounds on the step size parameter have not been previously derived. This creates obstacles in the total complexity analysis of such methods, because the oracle costs are typically decreasing in the step size parameter, and could be arbitrarily large as the step size parameter goes to 0. Thus, until now only the total iteration complexity of these methods has been analyzed. In this paper, we derive a lower bound on the step size parameter that holds with high probability for a large class of adaptive stochastic methods. We then use this lower bound to derive a framework for analyzing the expected and high probability total oracle complexity of any method in this class. Finally, we apply this framework to analyze the total sample complexity of two particular algorithms, STORM [BCMS19] and SASS [JSX21a], in the expected risk minimization problem.

## 1   Introduction

The widespread use of stochastic optimization algorithms for problems arising in machine learning and signal processing has made the stochastic gradient method and its variants overwhelmingly popular despite their theoretical and practical shortcomings. *Adaptive stochastic optimization* algorithms, on the other hand, borrow from decades of advances in deterministic optimization research, and offer new paths forward for stochastic optimization to be more effective and even more applicable. Adaptive algorithms can avoid many of the practical deficiencies of contemporary methods (such as the tremendous costs of tuning the step sizes of an algorithm for each individual application) while possessing strong convergence and worst-case complexity guarantees in surprisingly diverse settings.

Adaptive optimization algorithms have a long and successful history in deterministic optimization and include line search, trust-region methods, cubic regularized Newton methods, etc. All these methods have a common iterative framework, where at each iteration a candidate step is computed by the algorithm based on a local model of the objective function and a step size parameter that controls the length of this candidate step. The candidate step is then evaluated in terms of the decrease it achieves in the objective function, with respect to the decrease that it was expected to achieve based on the model. Whenever the decrease is sufficient, the step is accepted and the step size parameter may be increased to allow the next iteration to be more aggressive. If the decrease is not sufficient (or not achieved at all) then the step is rejected and a new step is computed using a smaller step size parameter. The model itself remains unchanged if it is known that a sufficiently small step will always succeed, as is true, for example with first- and second-order

---

[*]School of Operations Research and Information Engineering, Cornell University, Ithaca, NY, USA. Email: {bzj3,ks2375,mx229}@cornell.edu.

Taylor models of smooth functions. The analysis of such methods relies on the key property that the step size parameter is bounded away from zero and that once it is small enough the step is always accepted and the objective function gets reduced.

When stochastic oracles are used to estimate the objective function and its derivatives, the models no longer have the same property as the Taylor models; steps that decrease the model might not decrease the function, no matter how small the step size is. Thus all stochastic variants of these methods recompute the model at each iteration. The requirement on the stochastic model is then that it achieves a Taylor-like approximation with sufficiently high probability. This also means that steps may get rejected even if the step size parameter is small, simply because the stochastic oracles fail to deliver desired accuracy. Despite this difficulty, one can develop and analyze stochastic variants of adaptive optimization methods.

Recently developed stochastic variants of line search (which we will call step search, since unlike the deterministic version the search direction has to change on each iteration, thus the algorithm is not searching along a line) include [CS17, PS20, BCS19, JSX21a]. Stochastic trust-region methods have been analyzed in [BCMS19, GRVZ18, CBS22], and adaptive cubic regularized methods based on random models has been analyzed in [CS17, SX23]. For all these methods, bounds on iteration complexity have been derived, either in expectation or in high probability, under the assumption that stochastic oracles involved in estimating the objective function deliver sufficiently high accuracy with sufficiently high probability. While this probability is usually fixed, the accuracy requirement of the oracles is adaptive and depends on the step size parameter. Specifically, the smaller that parameter is, the more accurate the oracles need to be, in order to maintain the Taylor-like behavior of the model used by the algorithm. In most applications, having more accurate stochastic oracles implies a higher per-iteration cost. Furthermore, unlike the deterministic case, the step sizes in the stochastic case are not bounded away from zero due to possible oracle failures, and bounds on the step size parameter have not been previously derived. This creates significant difficulty in the analysis of the total oracle complexity of such methods because the oracle costs could be arbitrarily large as the step size parameter goes to zero.

In this paper, we derive a lower bound on the step size parameter for a general class of stochastic adaptive methods that encompasses all the algorithms in the preceding paragraph. This enables us to derive a bound on the total oracle complexity for any algorithm within this class and specific stochastic oracles arising, for example, from expected risk minimization. Our key contributions are as follows:

- Provide a high probability lower bound on the step size parameter for a wide class of stochastic adaptive methods using a coupling argument between the stochastic process generated by the algorithm and a one-sided random walk.

- Derive a framework for analyzing expected and high probability total oracle complexity bounds for this general class of stochastic adaptive methods.

- Apply these bounds to STORM (Stochastic Trust-region Optimization with Random Models) [BCMS19] and SASS (Stochastic Adaptive Step Search) [JSX21a] to derive their total sample complexity for expected risk minimization, and show they essentially match the complexity lower bound of first-order algorithms for stochastic non-convex optimization [ACD+19].

We note that there is a plethora of other algorithms in the literature that propose different ways to adaptively vary the step sizes in stochastic optimization. Some of these methods, such as [KB14, TMDQ16], have step size dynamics that are very different from the framework studied in this paper, hence the analysis of this paper is not needed or applicable. Some other adaptive methods, like [RVZ23, SHP18, RJM17], have step size dynamics that are more similar. However, since these papers do not provide iteration complexity bounds, we cannot obtain their sample complexity by directly applying our framework. Nonetheless, it would be interesting to explore whether our step size lower bound can be applied to the algorithms in these papers when a certain stopping time is specified.

We consider a continuous optimization problem of the form

$$\min_{x \in \mathbb{R}^m} \phi(x), \tag{1}$$

where $\phi$ is possibly non-convex, (twice-)continuously differentiable with Lipschitz continuous derivatives. Neither function values $\phi(x)$, nor gradients $\nabla \phi(x)$ are assumed to be directly computable. Instead, given

any $x \in \mathbb{R}^m$, it is assumed that stochastic estimates of $\phi(x)$, $\nabla\phi(x)$, and possibly $\nabla^2\phi(x)$ can be computed, and these estimates may possess different levels of *accuracy* and *reliability* depending on the particular setting of interest.

The adaptive stochastic algorithms that have been developed recently [CS17, PS20, BCS19, BCMS19, GRVZ18, JSX21a, CBS22, BXZ23, MWX23] use these stochastic estimates to compute an $\varepsilon$-optimal point $x_\varepsilon$, which means $\phi(x_\varepsilon) - \inf_x \phi(x) \le \varepsilon$ if $\phi$ is convex or $\|\nabla\phi(x_\varepsilon)\| \le \varepsilon$ if $\phi$ is non-convex. In the next section, we will introduce the general framework that encompasses these methods and discuss particular examples in more detail. In Section 2.6, we discuss the stochastic process generated by the algorithmic framework, including the stochastic step size parameter. In Section 3, we derive a lower bound on the step size parameter which holds in high probability. In Section 4, we use this lower bound to derive abstract expected and high probability total oracle complexity for any algorithm in this framework. In Section 5, we particularize these bounds to the specific examples of first-order STORM [BCMS19] and SASS [JSX21a] algorithms to bound their total sample complexity when applied to expected risk minimization. We conclude with some remarks in Section 6.

# 2 Algorithm framework and oracles

We now introduce and discuss an algorithmic framework for adaptive stochastic optimization in Algorithm 1. The framework is assumed to have access to stochastic oracles, that for any given point $x$ can generate random quantities $f(x, \xi_0) \approx \phi(x)$ (via zeroth-order oracle), $g(x, \xi_1) \approx \nabla\phi(x)$ (first-order oracle) and, possibly, $H(x, \xi_2) \approx \nabla^2\phi(x)$ (second-order oracle). After we introduce the algorithm we will discuss the general definition of an oracle that we use in this paper.

## 2.1 General Algorithm

At each iteration, given $x_k$, a stochastic local model $m_k(x_k + s) : \mathbb{R}^m \to \mathbb{R}$ of $\phi(x_k + s)$ is constructed using $g(x, \xi_1)$ and possibly $H(x, \xi_2)$. Using this model, a step $s_k(\alpha_k)$ is computed so that $m(x_k + s_k(\alpha_k))$ is a sufficient improvement over $m_k(x_k)$, where $\alpha_k$ is the step size parameter, which directly or indirectly controls the step length. The iterate $x_k$ and the trial point $x_k^+ = x_k + s_k(\alpha_k)$ are evaluated using the zeroth-order oracle values $f(x_k, \xi_{0,k})$ and $f(x_k^+, \xi_{0,k}^+)$. If these estimates suggest that sufficient improvement is attained, then the step is deemed *successful*, $x_k^+$ is accepted as the next iterate, and the parameter $\alpha_k$ is increased up to a multiplicative factor; otherwise, the step is deemed *unsuccessful*, the iterate does not change, and $\alpha_k$ is decreased by a multiplicative factor. Unlike in the deterministic case, new calls to all oracles are made at each iteration *even when the iterate does not change.*

For each method we give the form of $m_k(x_k + s)$, in terms of $g_k = g(x_k, \xi_{1,k})$ and $H_k = H(x_k, \xi_{2,k})$, $s_k(\alpha_k)$ and the sufficient reduction criterion in terms of $f_k^0 = f(x_k, \xi_{0,k})$ and $f_k^+ = f(x_k^+, \xi_{0,k}^+)$.

## 2.2 General Oracles

Typically in the optimization literature, an oracle is a computational procedure that provides the algorithm with some (estimate of) required information about the objective function. The oracle is endowed with some properties that the algorithm then utilizes. For example, an oracle may be assumed to produce $\nabla\phi(x)$ exactly – an assumption that a relevant optimization algorithm (and its analysis) will make use of and without which the algorithm may fail. Alternatively, an oracle may produce an inexact estimate of $\nabla\phi(x)$, in which case a relevant algorithm will operate under some specific knowledge about the properties of this estimator; e.g. that it is an unbiased estimator with variance that is bounded as a function of $\|\nabla\phi(x)\|$, or that it is a deterministic estimator with an error bounded by some known quantity. Algorithms and their analyses differ depending on the properties of the oracles they utilize. The algorithms analyzed in [CS17, PS20, BCS19, BCMS19, GRVZ18, JSX21a, SX23, CBS22, BXZ23, MWX23] all fit into the framework of Algorithm 1 but under a variety of specific assumptions on the oracles that generate function, gradient and and Hessian estimates. These assumptions are more complex than is typical in the literature, yet, they are shown to be applicable in many settings. For some extensive discussion on these properties, their comparison and specific examples we refer the reader, for example, to [CBS22]. Here we summarize the main ideas.

3

---

**Algorithm 1** **Algorithmic Framework for Adaptive Stochastic Optimization**

---

**0. Initialization**

Choose $\theta \in (0,1)$, $\gamma \in (0,1)$, $\alpha_{\max} \in (0,\infty)$, $x_0 \in \mathbb{R}^m$, and $\alpha_0 \in (0, \alpha_{\max}]$. Set $k \leftarrow 0$.

**1. Determine model and compute step**

Construct a stochastic model $m_k$ of $\phi$ at $x_k$ using $f(x_k, \xi_{0,k})$, $g(x_k, \xi_{1,k})$, and (optionally) $H(x_k, \xi_{2,k})$ from *probabilistic zeroth-, first-, and (optionally) second-order oracles*. Compute $s_k(\alpha_k)$ such that the model reduction $m_k(x_k) - m_k(x_k + s_k(\alpha_k)) \geq 0$ is sufficiently large.

**2. Check for sufficient reduction**

Set $x_k^+ \leftarrow x_k + s_k(\alpha_k)$ and compute $f(x_k^+, \xi_{0,k}^+)$ as a stochastic estimate of $\phi(x_k^+)$ using a *probabilistic zeroth-order oracle*. Check if $f(x_k, \xi_{0,k}) - f(x_k^+, \xi_{0,k}^+)$ is sufficiently large (e.g., relative to the model reduction $m_k(x_k) - m_k(x_k^+)$) using a condition parameterized by $\theta$.

**3. Successful iteration**

If sufficient reduction has been attained (along with other potential requirements), then set $x_{k+1} \leftarrow x_k^+$ and $\alpha_{k+1} \leftarrow \min\{\gamma^{-1}\alpha_k, \alpha_{\max}\}$.

**4. Unsuccessful iteration**

Otherwise, set $x_{k+1} \leftarrow x_k$ and $\alpha_{k+1} \leftarrow \gamma\alpha_k$.

**5. Next iteration**

Set $k \leftarrow k + 1$ and go to Step 1.

---

As presented above, the algorithmic framework described in Algorithm 1 has access to oracles that generate random quantities $f(x, \xi_0) \approx \phi(x)$ (zeroth-order oracle), $g(x, \xi_1) \approx \nabla\phi(x)$ (first-order oracle) and, possibly, $H(x, \xi_2) \approx \nabla^2\phi(x)$ (second-order oracle). Each of these oracles can have an input, an output and some intrinsic properties. The input is typically the current iterate $x_k$ and the step size parameter $\alpha_k$. These quantities, together with the intrinsic properties of the oracle, define the probability space (and thus the distribution) of the random variable $\xi$. For example, in the expected risk minimization setting, where the oracles are computed by averaging a minibatch of samples, $\xi_i$ (for $i \in \{0, 1, 2\}$) is the random minibatch used for the $i^{\text{th}}$-order oracle. The probability space and the distribution of $\xi_i$ may depend on $x$ and $\alpha$, since (as we will see in Section 5) the size of the minibatch may depend on these quantities. As another example, the first-order oracles can be computed using (randomized) finite differences [BCCS21]. In this case, the probability space and the distribution of $\xi_1$ may depend on the current iterate, the step size parameter and the randomness in function value estimates. More examples of stochastic oracles can include robust gradient estimation [Ans60, YCKB18], SPSA [Spa99], etc.

In summary, the oracles can be implemented in a variety of ways, depending on the application, while the algorithmic framework can be agnostic to how exactly the oracles are implemented. The algorithm operates under certain assumptions on the accuracy and reliability of the oracles. The cost of implementing an oracle to satisfy these requirements depends on the application, but it needs to be considered in an overall complexity analysis, which is what we address in this paper. The general oracles in this paper can be described at a high level as follows.

**Definition 1. Stochastic $j$th-order oracle.** *Given an input $x \in \mathbb{R}^n$ and the step size parameter $\alpha$, the oracle computes $\varphi_j(x, \xi_j)$, an estimate of the $j^{\text{th}}$-order derivative $\nabla^j\phi(x)$. In all the algorithms we consider, $\|\varphi_j(x, \xi_j) - \nabla^j\phi(x)\|$ is assumed to be bounded by some quantity (which will be a function of $\alpha$), with probability at least $1 - \delta_j$. Here, $\xi_j$ is a random variable defined on probability space $(\Omega_j, \mathcal{F}_j, P_j)$ whose distribution depends on the input $x$ and $\alpha$, and $\delta_j$ is intrinsic to the oracle. The cost of the oracle is also a function of $x$ and $\alpha$.*

In the next subsection, we describe how various methods fit into the general framework, and what specific requirements they have on the oracles.

## 2.3 Step search method

In the case of the step search (SS) methods in [CS17, BCS19, JSX21a] the particulars are as follows. Quantities $f_k^0$ $f_k^+$ and $g_k$ are random outputs of the stochastic oracles and $H_k$ is some positive definite matrix

(e.g., the identity).

- $m_k(x_k + s) = \phi(x_k) + g_k^T s + \frac{1}{2\alpha_k} s^T H_k s,$

- $s_k(\alpha_k) = -\alpha_k H_k^{-1} g_k$

- Sufficient reduction: $f_k^0 - f_k^+ \geq -\theta g_k^T s_k(\alpha_k) - r$

Note that although the true function value $\phi(x_k)$ appears in the definition of $m_k$, we do not need to know or query for it to run the algorithm because the function value is just a constant that makes no difference when we minimize the model. Here $\theta \in (0, 1)$, and $r$ is a small positive number that compensates for the noise in the function estimates. We will discuss the choice of $r$ after we introduce conditions on the oracle outputs $f(x, \xi_0)$ and $g(x, \xi_1)$.

- **SS.0** (Step search, zeroth-order oracle). Given a point $x$, the oracle computes a (random) function estimate $f(x, \xi_0)$ (where $\xi_0 = \xi_0(x)$ is the randomness of the oracle, which may depend on the current point $x$) such that
$$\mathbb{P}_{\xi_0}\left(|\phi(x) - f(x, \xi_0)| < \epsilon_f + t\right) \geq 1 - \delta_0(t),$$
  for some $\epsilon_f \geq 0$ and any $t > 0$.

- **SS.1** (Step search, first-order oracle). Given a point $x$ and the current *step size parameter* $\alpha$, the oracle computes a (random) gradient estimate $g(x, \xi_1)$ (where $\xi_1 = \xi_1(x, \alpha)$ is the randomness of the oracle, which may depend on $x$ and $\alpha$) such that
$$\mathbb{P}_{\xi_1}\left(\|g(x, \xi_1) - \nabla\phi(x)\| \leq \max\{\epsilon_g, \min\{\tau, \kappa\alpha\}\|g(x, \xi_1)\|\}\right) \geq 1 - \delta_1$$
  for some nonnegative constants $\epsilon_g$, $\kappa$, $\tau$ and $\delta_1$.

In [CS17], $\epsilon_f = 0$ and $\delta_0(t) \equiv 0$, which means that the zeroth-order oracle is exact, and $r = 0$. In [BCS19], $\epsilon_f > 0$ and $\delta_0(t) \equiv 0$, which means that the zeroth-order oracle has a bounded error with probability one, and $r = 2\epsilon_f$. In [JSX21a], $\epsilon_f > 0$ and $\delta_0(t) = e^{-\lambda t}$ for some $\lambda > 0$. This means there is no restriction on the error if it is less than $\epsilon_f$, and the tail of the error decays exponentially beyond $\epsilon_f$. In [JSX21a], $r > 2 \sup_x \mathbb{E}_{\xi_0}[|\phi(x) - f(x, \xi_0)|]$. In these works, $\theta$ is a parameter in $(0, 1)$ that can be chosen by the user running the algorithm.

In [CS17], $\epsilon_g = 0$ and $\delta_1 < \frac{1}{2}$. In [BCS19, JSX21a], $\epsilon_g > 0$ and $\delta_1$ is sufficiently small with a more complicated upper bound. In [CS17, BCS19], $\alpha_{\max}$ is finite, thus $\tau$ is $\kappa\alpha_{\max}$. In [JSX21a], $\alpha_{\max}$ is infinity, and $\tau$ is simply assumed to be some constant intrinsic to the oracle.

## 2.4 Trust-region method

Stochastic trust-region (TR) methods that fall into the framework of Algorithm 1 have been developed and analyzed in [GRVZ18, BCMS19, CBS22]. In the case of TR algorithms, $f_k^0$, $f_k^+$, $g_k$, and (possibly) $H_k$ are random outputs of the stochastic oracles, and

- $m_k(x_k + s) = \phi(x_k) + g_k^T s + \frac{1}{2} s^T H_k s,$

- $s_k(\alpha_k) = \arg\min_{s:\, \|s\| \leq \alpha_k} m_k(x_k + s)$

- Sufficient reduction: $\frac{f_k^0 - f_k^+ + r}{m_k(x_k) - m_k(x_k + s_k(\alpha_k))} \geq \theta$

- Additional requirement for a successful iteration: $\|g_k\| \geq \theta_2 \alpha_k$, for some $\theta_2 > 0$.

The requirements for the oracles are as follows. In the case of first-order analysis in [GRVZ18, BCMS19, CBS22], the following first-order oracle is assumed to be available.

- **TR1.1** (First-order trust-region, first-order oracle). Given a point $x$ and the current *trust-region radius* $\alpha$, the oracle computes a gradient estimate $g(x, \xi_1)$ (where $\xi_1 = \xi_1(x, \alpha)$ is the randomness of the oracle, which can depend on $x$ and $\alpha$) such that

$$\mathbb{P}_{\xi_1}\left(\|g(x, \xi_1) - \nabla\phi(x)\| \leq \epsilon_g + \kappa_{eg}\alpha\right) \geq 1 - \delta_1.$$

  Here, $\kappa_{eg}$ and $\delta_1$ are nonnegative constants.

In the second-order analysis, the following first- and second-order oracles are used:

- **TR2.1** (Second-order trust-region, first-order oracle). Given a point $x$ and the current *trust-region radius* $\alpha$, the oracle computes a gradient estimate $g(x, \xi_1)$ (where $\xi_1 = \xi_1(x, \alpha)$ is the randomness of the oracle, which can depend on $x$ and $\alpha$) such that

$$\mathbb{P}_{\xi_1}\left(\|g(x, \xi_1) - \nabla\phi(x)\| \leq \epsilon_g + \kappa_{eg}\alpha^2\right) \geq 1 - \delta_1.$$

  Here, $\kappa_{eg}$ and $\delta_1$ are nonnegative constants.

- **TR2.2** (Second-order trust-region, second-order oracle). Given a point $x$ and the current *trust-region radius* $\alpha$, the oracle computes a Hessian estimate $H(x, \xi_2)$ (where $\xi_2 = \xi_2(x, \alpha)$ is the randomness of the oracle, which can depend on $x$ and $\alpha$) such that

$$\mathbb{P}_{\xi_2}\left(\|H(x, \xi_2) - \nabla\phi(x)\| \leq \epsilon_h + \kappa_{eh}\alpha\right) \geq 1 - \delta_2.$$

  Here, $\kappa_{eh}$ and $\delta_2$ are nonnegative constants.

$\epsilon_h$ and $\epsilon_g$ are assumed to equal 0 in [GRVZ18, BCMS19] but are allowed to be positive in [CBS22].

In terms of the zeroth-order oracles, the three works make different assumptions. Specifically, in [GRVZ18], as in [CS17], the zeroth-order oracle is assumed to be exact. In [CBS22] the zeroth-order oracle is the same as in [JSX21a], and $r > 2\epsilon_f + \frac{2}{\lambda}\log 4$. For the first-order analysis in [BCMS19], however, the zeroth-order oracle is as follows (and $r = 0$).

- **TR1.0** (First-order trust-region, zeroth-order oracle). Given a point $x$ and the current *trust-region radius* $\alpha$, the oracle computes a function estimate $f(x, \xi_0)$ (where $\xi_0 = \xi_0(x, \alpha)$ is the randomness of the oracle, which can depend on $x$ and $\alpha$) such that

$$\mathbb{P}_{\xi_0}\left(|f(x, \xi_0) - \phi(x)| \leq \kappa_{ef}\alpha^2\right) \geq 1 - \delta_0,$$

  where $\kappa_{ef}$ and $\delta_0$ are some nonnegative constants.

For the second-order analysis in [BCMS19], the zeroth-order oracle requirements are tighter.

- **TR2.0** (Second-order trust-region, zeroth-order oracle). Given a point $x$ and the current *trust-region radius* $\alpha$, the oracle computes a function estimate $f(x, \xi_0)$ (where $\xi_0 = \xi_0(x, \alpha)$ is the randomness of the oracle, which can depend on $x$ and $\alpha$) such that

$$\mathbb{P}_{\xi_0}\left(|f(x, \xi_0) - \phi(x)| \leq \kappa_{ef1}\alpha^3\right) \geq 1 - \delta_0$$

  and

$$\mathbb{E}_{\xi_0}\left[|f(x, \xi_0) - \phi(x)|\right] \leq \kappa_{ef2}\alpha^3$$

  where $\kappa_{ef1}$, $\kappa_{ef2}$ and $\delta_0$ are some nonnegative constants.

## 2.5 Cubicly regularized Newton method

The cubicly regularized (CR) Newton methods in [CS17, SX23] also fit the framework of Algorithm 1 with

- $m_k(x_k + s) = \phi(x_k) + g_k^T s + \frac{1}{2}s^T H_k s + \frac{1}{3\alpha_k}\|s\|^3$,

- $s_k(\alpha_k) = \arg\min_s m_k(x_k + s),$ $\hspace{4cm}$ (2)

- Sufficient reduction: $\frac{f_k^0 - f_k^+ + r}{m_k(x_k) - m_k(x_k + s_k(\alpha_k))} \geq \theta$.

In [CS17], the zeroth-order oracle is assumed to be exact, that is $f_k^0 = \phi(x_k)$ and $f_k^+ = \phi(x_k + s_k(\alpha_k))$, and $r = 0$. In [SX23] the zeroth-order oracle and the choice of $r$ are the same as in [JSX21a]. In [CS17] and [SX23], a version of the following first- and second-order oracles are used. For specific implementation details, please refer to [CS17, SX23].

- **CR.1** (Cubicly regularized Newton, first-order oracle). Given a point $x$ and the current parameter $\alpha$, the oracle computes a gradient estimate $g(x, \xi_1)$ (where $\xi_1 = \xi_1(x, \alpha)$ is the randomness of the oracle, which can depend on $x$ and $\alpha$) such that

$$\mathbb{P}_{\xi_1}\left(\|\nabla\phi(x) - g(x, \xi_1))\| \leq \kappa_{eg} \max\left\{\alpha, \|s\|^2\right\}\right) \geq 1 - \delta_1,$$

where $\kappa_{eg}$ and $\delta_1$ are nonnegative constants, and $s$ is defined in (2).

- **CR.2** (Cubicly regularized Newton, second-order oracle). Given a point $x$, and the current parameter $\alpha$, the oracle computes a Hessian estimate $H(x, \xi_2)$ (where $\xi_2 = \xi_2(x, \alpha)$ is the randomness of the oracle, which can depend on $x$ and $\alpha$) such that

$$\mathbb{P}_{\xi_2}\left(\|(\nabla^2\phi(x) - H(x, \xi_2))s\| \leq \kappa_{eh} \max\left\{\alpha, \|s\|^2\right\}\right) \geq 1 - \delta_2,$$

where $\kappa_{eh}$ and $\delta_2$ are nononegative constants, and $s$ is defined in (2).

It is apparent that all of the algorithms that we discussed above rely on oracles whose accuracy requirements change adaptively with $\alpha$. It is also clear that for many settings, a higher accuracy requirement leads to a higher oracle complexity. For example, if a stochastic oracle is delivered via sample averaging, then more samples are needed to provide a higher accuracy. Therefore, to bound the total oracle complexity of the algorithm, we need to bound the accuracy requirement over the iterations, or equivalently provide a lower bound for the parameter $\alpha$.

## 2.6 Notions of the stochastic process

When applied to problem (1), Algorithm 1 generates a stochastic process (with respect to the randomness underlying the stochastic oracles). Specifically, let $(X_k)_{k\geq 0}$ be the random iterates with realizations $x_k$, let $(G_k)_{k\geq 0}$ be the gradient estimates with realizations $g_k$, and let $(\mathcal{A}_k)_{k\geq 0}$ be the step size parameter values with realizations $\alpha_k$. The prior works that analyze different methods belonging to the framework of Algorithm 1 define this stochastic process rigorously, with appropriate filtrations. Here for brevity, we will omit those details, as we do not use them in the analysis. We now define a stopping time for the process.

**Definition 1** (Stopping time). *For $\varepsilon > 0$, let $T_\varepsilon$ be the first time such that a specified optimality condition is satisfied. For all the settings considered in this paper, $T_\varepsilon = \min\{k : \|\nabla\phi(x_k)\| \leq \varepsilon\}$ if $\phi$ is non-convex, and $T_\varepsilon = \min\{k : \phi(x_k) - \inf_x \phi(x) \leq \varepsilon\}$ if $\phi$ is strongly convex. We will refer to $T_\varepsilon$ as the stopping time of the algorithm.*

The following property is crucial in the analysis of algorithms in the framework of Algorithm 1.

**Assumption 1** (Properties of the stochastic process generated by the adaptive stochastic algorithm). *The random sequence of parameters $\mathcal{A}_k$ generated by the algorithm satisfies the following:*

(i) *For all $k$, $\mathcal{A}_k \in \{\gamma\mathcal{A}_{k-1}, \min\{\alpha_{\max}, \gamma^{-1}\mathcal{A}_{k-1}\}\}$, and*

(ii) *There exist constants $\bar{\alpha} > 0$, and $p > \frac{1}{2}$, such that for all iterations $k$,*

$$\mathbb{P}(\mathcal{A}_{k+1} = \gamma^{-1}\mathcal{A}_k \mid \mathcal{F}_k, k < T_\varepsilon, \mathcal{A}_k \leq \bar{\alpha}) \geq p.$$

*Here, $\mathcal{F}_k$ denotes the filtration generated by the algorithm up to iteration $k$.*

The algorithms in [CS17, PS20, BCS19, BCMS19, GRVZ18, JSX21a, CBS22] all satisfy Assumption 1, under appropriate lower bounds on the oracle probabilities $\delta_0, \delta_1$ (and $\delta_2$). Also, without loss of generality, we will assume that $\alpha_0 \geq \bar{\alpha}$, since if $\alpha_0 < \bar{\alpha}$ we can simply define $\bar{\alpha}$ to be $\alpha_0$, as $\alpha_0$ is a constant. In the next section, under Assumption 1, we derive a high probability lower bound on $\alpha_k$ as a function of the number of iterations $n$, $\bar{\alpha}$, $p$, and $\gamma$.

Throughout the remainder of this paper, we will use $q$ to denote $1 - p$.

# 3 High probability lower bound for the step size parameter

The following theorem provides a high probability lower bound for $\alpha_k$.

**Theorem 1.** *Suppose Assumption 1 holds for Algorithm 1. For any positive integer $n$, any $\omega > 0$, with probability at least $1 - n^{-\omega} - cn^{-(1+\omega)}$, we have*
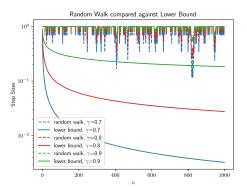
$$\text{either } T_\varepsilon < n \quad \text{or} \quad \min_{1 \leq k \leq n} \alpha_k \geq \alpha^*(n) := \bar{\alpha}\gamma\gamma^{(1+\omega)\log_{1/2q} n} = \bar{\alpha}\gamma n^{-(1+\omega)\log_{1/2q} 1/\gamma},$$

*where $c = \frac{2\sqrt{pq}}{(1-2\sqrt{pq})^2}$ and $q = 1 - p$.*

The proof of this theorem involves two steps. First, in Section 3.1, we show that for $n < T_\varepsilon$, the sequence of step size parameters $\mathcal{A}_k$ generated by the algorithm can be coupled with a random walk on the non-negative integers. This reduces the problem to that of bounding the maximum value of a one-sided random walk in the first $n$ steps. We then derive a high probability upper bound on this maximum value in Section 3.2.

Before moving to its proof, we illustrate the theorem using some plots and comment on some implications of the theorem.

**Illustration of Theorem 1.** Figure 1a illustrates the high probability bound provided by Theorem 1. The solid curves depict the lower bounds given by the theorem for $\bar{\alpha} = 1, \omega = 1$, $p = 0.8$, and for varying values of $\gamma$. In comparison, the dotted lines correspond to one-sided random walks $\mathcal{Z}_k$ that start at $\bar{\alpha} = 1$. At each step, $\mathcal{Z}_{k+1} = \gamma\mathcal{Z}_k$ with probability $1 - p$, and $\mathcal{Z}_{k+1} = \min\{1, \gamma^{-1}\mathcal{Z}_k\}$ with probability $p$. The proof of Theorem 1 shown later implies that there is a coupling between the sequence of parameters $\mathcal{A}_k$ generated by the algorithm and $\mathcal{Z}_k$, such that $\mathcal{A}_k \geq \mathcal{Z}_k$, in other words, the sequence of parameters $\mathcal{A}_k$ generated by the algorithm *stochastically dominates $\mathcal{Z}_k$*.



(a) Comparing the random walk trajectory with the theoretical lower bound for various values of $\gamma$.

(b) Comparing the *minimum value attained so far by the random walk curves in Figure 1a* with the theoretical lower bounds (the same as in Figure 1a).

Figure 1: Illustration of Theorem 1.

**Remarks on Theorem 1.**

1. For fixed $n, \gamma$, and $\bar{\alpha}$, the lower bound is a function of $p$. It increases as $p$ increases. Specifically, the exponent of $n$ changes with $p$, and the exponent goes to 0 as $p$ goes to 1. Hence as $p$ goes to 1, this lower bound simplifies to $\bar{\alpha}\gamma$, which matches the lower bound in the deterministic case.

2. When $p$ is close to 1 (i.e. when the stochastic oracles are highly reliable), this lower bound decreases slowly as a function of $n$, since the exponent of $n$ is close to 0. Alternatively, when the stochastic oracles are not highly reliable, increasing the value of $\gamma$ allows the algorithm to maintain a slow decrease of the step size.

3. Enlarging $\gamma$ as $p$ decreases makes intuitive sense for the algorithm. When $p$ is large, an unsuccessful step is more likely to be caused by the step size being too large rather than the failure of the oracles to deliver the desired accuracy. On the other hand, when $p$ is small, unsuccessful iterations are likely to occur even when the step size parameter is already small. Thus in the latter case, larger $\gamma$ values help avoid an erroneous rapid decrease of the step size parameter.

4. If we choose $\gamma = \left(\frac{1}{2q}\right)^{-\frac{1}{4}}$ and $\omega = 1$, then the minimum step size is lower bounded by $\bar{\alpha}\gamma n^{-\frac{1}{2}}$ with high probability. This coincides with the typical choice of the step size decay schemes for the stochastic gradient method applied to non-convex functions.

The theorem implies that we can even bound $\alpha$ by a *constant* times $\bar{\alpha}$ with high probability, provided we set $\gamma$ as a function of $n$.

**Corollary 1.** *Let Assumption 1 hold for Algorithm 1, then for any positive integer $n$, any $\omega > 0$, and any $\beta < \frac{1}{2}$, if*

$$\gamma \geq \max\left\{\frac{1}{2}, \left(\frac{1}{2q}\right)^{\frac{\log(2\beta)}{(1+\omega)\log n}}\right\},$$

*then with probability at least $1 - n^{-\omega} - cn^{-(1+\omega)}$, where $c = \frac{2\sqrt{pq}}{(1-2\sqrt{pq})^2}$, we have*

$$\text{either } T_\varepsilon < n \text{ or } \min_{1 \leq k \leq n} \alpha_k \geq \beta\bar{\alpha}.$$

*Proof.* This follows from Theorem 1 by substituting in the specified value of $\gamma$. □

In the remainder of this section, we prove Theorem 1 in two steps.

## 3.1 Step 1: reduction to random walk

We will use a coupling argument to obtain the reduction to a random walk.

Let $\{\mathcal{A}_k\}_{k=0}^\infty$ denote the random sequence of parameter values (whose realization is $\{\alpha_k\}_{k=0}^\infty$), for Algorithm 1. Let us assume, WLOG, that $\mathcal{A}_0 = \gamma^j \bar{\alpha}$, for some integer $j \leq 0$. (Recall here that $0 < \gamma < 1$.) Then we observe that $\mathcal{A}_k = \gamma^{\bar{Y}_k}\bar{\alpha}$, where $\{\bar{Y}_k\}_{k=0}^\infty$ is a random sequence of integers, with $\bar{Y}_0 = j \leq 0$, which increases by one on every unsuccessful step, and decreases by one on every successful step. Moreover, by Assumption 1, whenever $k < T_\varepsilon$ and $\bar{Y}_k \geq 0$, the probability that it decreases by one is at least $p$. Define $Y_k$ as follows:

$$Y_k = \bar{Y}_k \text{ if } k \leq T_\varepsilon, \qquad Y_k = \begin{cases} Y_{k-1} - 1 & \text{w.p. } p \\ Y_{k-1} + 1 & \text{w.p. } 1 - p \end{cases} \text{ if } k > T_\varepsilon. \tag{3}$$

In other words, $Y_k$ follows the algorithm until $T_\varepsilon$, and then behaves like a random walk with downward drift $p$ after $T_\varepsilon$. We now couple $\{Y_k\}_{k=0}^\infty$ with a random walk $\{Z_k\}_{k=0}^\infty$ which stochastically dominates $Y_k$.

Consider the following one-sided random walk $\{Z_k\}_{k=0}^\infty$, defined on the non-negative integers.

$$Z_0 = 0, \qquad Z_{k+1} = \begin{cases} Z_k + 1, & \text{w.p. } 1 - p, \\ Z_k - 1, & \text{w.p. } p, \text{ if } Z_k \geq 1, \\ 0, & \text{w.p. } p, \text{ if } Z_k = 0. \end{cases} \tag{4}$$

**Lemma 1.** *There exists a coupling between $Z_k$ and $Y_k$, where $Z_k$ stochastically dominates $Y_k$.*

*Proof of Lemma 1.* Initially, $Z_0 = 0$ and $Y_0 \leq 0$. For each $k$, we show how to update $Z_k$ to $Z_{k+1}$ according to how $Y_k$ changes to $Y_{k+1}$. We consider two cases depending on whether $k < T_\varepsilon$ or $k \geq T_\varepsilon$.

Case 1: $k < T_\varepsilon$. If $Y_k \leq -1$, we update $Z_{k+1}$ from $Z_k$ according to Equation (4), independently of how $Y_k$ changes to $Y_{k+1}$. If $Y_k \geq 0$, then we first check if $Y_k$ increased or decreased. Let $p'$ be the probability that $Y_{k+1} = Y_k - 1$ on this sample path. Since $Y_k \geq 0$, we know by Assumption 1 that $p' \geq p$. Now, if $Y_{k+1} = Y_k + 1$, then we set $Z_{k+1} = Z_k + 1$. On the other hand, if $Y_{k+1} = Y_k - 1$, then we set $Z_{k+1} = Z_k + 1$

9

with probability $1 - \frac{p}{p'}$, and $Z_{k+1} = \max\{Z_k - 1, 0\}$ with probability $\frac{p}{p'}$. Note that these probabilities are well-defined because $p' \geq p$.

Case 2: $k \geq T_\varepsilon$. If $Y_{k+1} = Y_k + 1$, then set $Z_{k+1} = Z_k + 1$. Otherwise, if $Y_{k+1} = Y_k - 1$, then set $Z_{k+1} = \max\{Z_k - 1, 0\}$.

Observe that under this coupling, $Z_k \geq Y_k$ on every sample path. Moreover, $\{Z_k\}$ and $\{Y_k\}$ have the correct marginal distributions. For $Y_k$, this is easy to see, since it evolves according to its true distribution and we are constructing $Z_k$ from it. For $Z_k$, on any step with $k \geq T_\varepsilon$, $Z_{k+1}$ evolves from $Z_k$ correctly according to Equation (4) by construction. On a step with $k < T_\varepsilon$ there are two cases: 1) $Y_k \leq -1$, and 2) $Y_k \geq 0$. In the first case, the update from $Z_k$ to $Z_{k+1}$ clearly follows Equation (4). This is also true in the second case, since there the probability that $Z_k$ increases is $(1 - p') + p'(1 - \frac{p}{p'}) = 1 - p$.

To summarize, we have exhibited a coupling between $\{Z_k\}$ and $\{Y_k\}$, under which $Z_k \geq Y_k$ on any sample path. $\qquad\square$

## 3.2 Step 2: upper-bounding the maximum value of the random walk

We now derive a high probability upper bound on the maximum value reached by the random walk.

**Definition 2.** *Let $N(\ell, n)$ be the random variable that denotes the number of times $Z_k = \ell$ in the first $n$ steps of the random walk.*

By definition of $N(\ell, n)$, we have $N(\ell, n) > 0$ if and only if state $\ell$ is visited in the first $n$ steps of the random walk. The next proposition upper bounds the probability that $N(\ell, n) > 0$.

**Proposition 1.** *Let $q = 1 - p$. We have*

$$\mathbb{P}(N(\ell, n) > 0) \leq (n - \ell + 1)\frac{1 - (q/p)}{1 - (q/p)^{\ell+1}}\left(\frac{q}{p}\right)^\ell + \frac{2\sqrt{pq}}{(1 - 2\sqrt{pq})^2}(2q)^\ell.$$

*Proof.* First, observe that $\mathbb{P}(N(\ell, n) > 0)$ remains unchanged if we change the state space from $\{0, 1, 2, \ldots\}$ to $\{0, 1, 2, \ldots, \ell\}$ and modify the walk to hold in state $\ell$ with probability $q$ (instead of moving from $\ell$ to $\ell + 1$ with that probability). This defines a Markov chain on $\{0, 1, 2, \ldots, \ell\}$, and let $P$ be its transition matrix. Noting that $P_{0,\ell}^m$ is the probability that the Markov chain is in state $\ell$ at time $m$, we see that

$$\mathbb{P}(N(\ell, n) > 0) \leq \sum_{m=\ell}^n P_{0,\ell}^m. \tag{5}$$

The matrix $P$ is explicitly diagonalized in [Fel68] (Section XVI.3). By (3.16) in that section,

$$P_{0,\ell}^m = \frac{1 - (q/p)}{1 - (q/p)^{\ell+1}}\left(\frac{q}{p}\right)^\ell - \frac{2q}{\ell+1}\left(\frac{q}{p}\right)^{\frac{\ell-1}{2}}\sum_{r=1}^\ell \frac{\left[\sin\frac{\pi r}{\ell+1}\right]\left[\sin\frac{\pi r\ell}{\ell+1}\right]\left[2\sqrt{pq}\cos\frac{\pi r}{\ell+1}\right]^m}{1 - 2\sqrt{pq}\cos\frac{\pi r}{\ell+1}}. \tag{6}$$

The absolute value of the sum appearing in (6) can of course be bounded above by

$$\sum_{r=1}^\ell \frac{(2\sqrt{pq})^m}{1 - 2\sqrt{pq}} = \ell\frac{(2\sqrt{pq})^m}{1 - 2\sqrt{pq}}$$

and this readily yields

$$P_{0,\ell}^m \leq \frac{1 - (q/p)}{1 - (q/p)^{\ell+1}}\left(\frac{q}{p}\right)^\ell + \frac{2p}{1 - 2\sqrt{pq}}\left(\frac{q}{p}\right)^{\frac{\ell+1}{2}}(2\sqrt{pq})^m. \tag{7}$$

Summing (7) over $m = \ell, \ldots, n$ and using (5), we obtain the bound on $\mathbb{P}(N(\ell, n) > 0)$ claimed in the proposition. $\qquad\square$

**Remark 1.** *The bound for Proposition 1 is essentially tight, as the decay of $\mathbb{P}(N(\ell, n) > 0)$ is not faster than geometric; $q^\ell$ is a lower bound.*

With the above proposition at hand, Theorem 1 is proved by choosing an appropriate level $\ell$, for which $\mathbb{P}(N(\ell, n) = 0)$ is high.

*Proof of Theorem 1.* Let $q = 1 - p$. By Proposition 1, we have:

$$\mathbb{P}(N(\ell, n) > 0) \leq (n - \ell + 1)\frac{1 - (q/p)}{1 - (q/p)^{\ell+1}}\left(\frac{q}{p}\right)^\ell + \frac{2\sqrt{pq}}{(1 - 2\sqrt{pq})^2}(2q)^\ell.$$

In other words,

$$\mathbb{P}(N(\ell, n) = 0) \geq 1 - (n - \ell + 1)\frac{1 - (q/p)}{1 - (q/p)^{\ell+1}}\left(\frac{q}{p}\right)^\ell - \frac{2\sqrt{pq}}{(1 - 2\sqrt{pq})^2}(2q)^\ell$$

$$\geq 1 - n\left(\frac{q}{p}\right)^\ell - \frac{2\sqrt{pq}}{(1 - 2\sqrt{pq})^2}(2q)^\ell.$$

Let $a > 0$ be a parameter to be set later, and take $\ell = \lceil a \log(n) \rceil$. Then, the above inequality implies:

$$\mathbb{P}(N(\ell, n) = 0) \geq 1 - n\left(\frac{q}{p}\right)^{a \log(n)} - \frac{2\sqrt{pq}}{(1 - 2\sqrt{pq})^2}(2q)^{a \log(n)}$$

$$= 1 - n^{1 - a \log(\frac{p}{q})} - \frac{2\sqrt{pq}}{(1 - 2\sqrt{pq})^2}n^{-a \log(1/(2q))}.$$

In going from the first line to the second, we used the following fact about logarithms: $x^{a \log n} = n^{a \log x} = n^{-a \log(1/x)}$. Let $a = \frac{1+\omega}{\log(1/2q)}$. Then, $a \log(\frac{p}{q}) \geq a \log(\frac{1}{2q}) = 1 + \omega$, so

$$\mathbb{P}(N(\ell, n) = 0) \geq 1 - n^{-\omega} - \frac{2\sqrt{pq}}{(1 - 2\sqrt{pq})^2}n^{-(1+\omega)}.$$

In other words, with probability at least $1 - n^{-\omega} - cn^{-(1+\omega)}$ with $c = \frac{2\sqrt{pq}}{(1-2\sqrt{pq})^2}$, the random walk will remain below $\ell = \left\lceil (1 + \omega) \log_{1/2q} n \right\rceil$ in the first $n$ steps. By construction of the coupling, with the above probability, we know the $\alpha$ parameter in the algorithm either remains above $\gamma^{\lceil a \log n \rceil}\bar{\alpha} = \gamma^{\left\lceil (1+\omega) \log_{1/2q} n \right\rceil}\bar{\alpha}$ throughout the first $n$ steps, or the algorithm has reached its stopping time in $n$ steps.

$\square$

It is natural to ask if the theory extends to the case where the factors for increasing and decreasing the step size are different. The main challenge in this case is that the stochastic process of the step sizes is now modeled by a random walk whose step length going up is different than the step length going down. If it turns out that the hitting probability $\mathbb{P}(N(\ell, n) > 0)$ can be bounded for the more general one-sided random walks, we believe the theory in our paper can be extended similarly. We will leave it as a subject for future research.

## 4  Expected and high probability total oracle complexity

We now use the tools derived in the previous section to obtain abstract expected and high probability upper bounds on the total oracle complexity of Algorithm 1. In Section 5, we will derive concrete bounds for the total oracle complexity of two specific algorithms (STORM and SASS), and the specific oracles arising in expected risk minimization.

The cost of an oracle call may depend on the step size parameter $\alpha$ and the probability parameter $1 - \delta$, thus we denote the cost by $\mathsf{oc}(\alpha, 1 - \delta)$. We will use $\mathsf{oc}(\alpha)$ in the paper to simplify the notation because for all algorithms in the class, $\delta$ can be treated as a constant. Moreover, the cost of an oracle call is a non-increasing function of $\alpha$ for all algorithms developed so far that fit into the framework.

**Assumption 2.** $\mathsf{oc}(\alpha)$ *is non-increasing in* $\alpha$.

**Definition 3** (Total Oracle Complexity). *For a positive integer $n$, let $\mathsf{TOC}(n)$ be the random variable which denotes the total oracle complexity of running the algorithm for $n$ iterations. In other words,*

$$\mathsf{TOC}(n) = \sum_{k=1}^{n} \mathsf{oc}(\mathcal{A}_k).$$

## 4.1 Abstract expected total oracle complexity

We now proceed to bound $\mathsf{TOC}(\min\{T_\epsilon, n\})$ in expectation, where $n$ is an arbitrary positive integer.

**Theorem 2.** *Let Assumption 1 and Assumption 2 hold in Algorithm 1. For any positive integer $n$, we have*

$$\mathbb{E}[\mathsf{TOC}(\min\{T_\epsilon, n\})] \leq n \sum_{\ell=1}^{n} \min\left\{1, n\left(\frac{q}{p}\right)^\ell + \frac{2\sqrt{pq}}{(1-2\sqrt{pq})^2}(2q)^\ell\right\} \cdot \mathsf{oc}(\bar{\alpha}\gamma^\ell) + n\,\mathsf{oc}(\bar{\alpha}).$$

*Proof.* First, observe that if the $\alpha_k$ parameters are all above some value $\alpha^*$ in the first $n$ steps, then by Assumption 2, $\mathsf{TOC}(n) \leq n \cdot \mathsf{oc}(\alpha^*)$. Therefore, for any integer $\ell \geq 0$, we have

$$\mathbb{P}(\mathsf{TOC}(\min\{T_\epsilon, n\}) > n \cdot \mathsf{oc}(\bar{\alpha}\gamma^\ell)) \leq \mathbb{P}(N(\ell+1, n) > 0). \tag{8}$$

By Proposition 1,

$$\mathbb{P}(N(\ell+1, n) > 0) \leq n\left(\frac{q}{p}\right)^{\ell+1} + \frac{2\sqrt{pq}}{(1-2\sqrt{pq})^2}(2q)^{\ell+1}.$$

This implies

$$\mathbb{P}(\mathsf{TOC}(\min\{T_\epsilon, n\}) > n \cdot \mathsf{oc}(\bar{\alpha}\gamma^\ell)) \leq \min\left\{1, n\left(\frac{q}{p}\right)^{\ell+1} + \frac{2\sqrt{pq}}{(1-2\sqrt{pq})^2}(2q)^{\ell+1}\right\}.$$

By the definition of expectation,

$$
\begin{aligned}
&\mathbb{E}[\mathsf{TOC}(\min\{T_\epsilon, n\})] \\
&= \sum_{i=0}^{n\cdot\mathsf{oc}(\bar{\alpha}\gamma^n)} i \cdot \mathbb{P}(\mathsf{TOC}(\min\{T_\epsilon, n\}) = i) \\
&\leq \sum_{\ell=0}^{n-1} \mathbb{P}\left(\mathsf{TOC}(\min\{T_\epsilon, n\}) \in (n \cdot \mathsf{oc}(\bar{\alpha}\gamma^\ell), n \cdot \mathsf{oc}(\bar{\alpha}\gamma^{\ell+1})]\right) \cdot n\,\mathsf{oc}(\bar{\alpha}\gamma^{\ell+1}) \\
&\qquad + \mathbb{P}\left(\mathsf{TOC}(\min\{T_\epsilon, n\}) \in [0, n\,\mathsf{oc}(\bar{\alpha})]\right) \cdot n\,\mathsf{oc}(\bar{\alpha}) \\
&\leq \sum_{\ell=0}^{n-1} \mathbb{P}\left(\mathsf{TOC}(\min\{T_\epsilon, n\}) > n\,\mathsf{oc}(\bar{\alpha}\gamma^\ell)\right) \cdot n\,\mathsf{oc}(\bar{\alpha}\gamma^{\ell+1}) + n\,\mathsf{oc}(\bar{\alpha}) \\
&\leq \sum_{\ell=0}^{n-1} \min\left\{1, n\left(\frac{q}{p}\right)^{\ell+1} + \frac{2\sqrt{pq}}{(1-2\sqrt{pq})^2}(2q)^{\ell+1}\right\} \cdot n\,\mathsf{oc}(\bar{\alpha}\gamma^{\ell+1}) + n\,\mathsf{oc}(\bar{\alpha}).
\end{aligned}
$$

$\square$

## 4.2 Abstract high probability total oracle complexity

We now proceed to bound $\mathsf{TOC}(T_\varepsilon)$ in high probability, using Theorem 1.

12

**Theorem 3.** *Let Assumption 1 and Assumption 2 hold in Algorithm 1. For any $\omega > 0$ and positive integer $n$, with probability at least $1 - \mathbb{P}(T_\varepsilon > n) - n^{-\omega} - cn^{-(1+\omega)}$,*

$$\mathsf{TOC}(T_\varepsilon) \leq n \cdot \mathsf{oc}(\alpha^*(n)),$$

*where $\alpha^*(n) = \bar{\alpha}\gamma n^{-(1+\omega)\log_{1/2q} 1/\gamma}$ , and $c$ is as defined in Theorem 1.*

*If $\gamma$ is chosen to be at least $\max\left\{\frac{1}{2}, \left(\frac{1}{2q}\right)^{\frac{\log(2\beta)}{(1+\omega)\log n}}\right\}$ for some $\beta < \frac{1}{2}$, then $\alpha^*(n) \geq \beta\bar{\alpha}$, thus with probability at least $1 - \mathbb{P}(T_\varepsilon > n) - n^{-\omega} - cn^{-(1+\omega)}$,*

$$\mathsf{TOC}(T_\varepsilon) \leq n \cdot \mathsf{oc}(\beta\bar{\alpha}).$$

*Proof.* Let $\mathsf{TOC}_{\mathsf{rw}}(n)$ be the total oracle complexity of the first $n$ iterations with the corresponding sequence of parameters $\alpha_k$ induced by the one-sided random walk (that is, the sequence defined by $\alpha_k = \bar{\alpha}\gamma^{Z_k}$, where $Z_k$ is defined in Section 3.1). In other words,

$$\mathsf{TOC}_{\mathsf{rw}}(n) = \sum_{k=1}^{n} \mathsf{oc}(\bar{\alpha}\gamma^{Z_k}).$$

With probability $1 - \mathbb{P}(T_\varepsilon > n)$, we have $T_\varepsilon \leq n$, which implies

$$\mathsf{TOC}(T_\varepsilon) \leq \mathsf{TOC}_{\mathsf{rw}}(T_\varepsilon) \leq \mathsf{TOC}_{\mathsf{rw}}(n).$$

Here, the first inequality is by Lemma 1 and Assumption 2, and the second inequality is by $T_\varepsilon \leq n$.

The same arguments used in the proof of Theorem 1 show that with probability at least $1 - n^{-\omega} - cn^{-(1+\omega)}$, we have $\min_{1 \leq k \leq n} \bar{\alpha}\gamma^{Z_k} \geq \alpha^*(n)$. Thus, with at least this probability, $\mathsf{TOC}_{\mathsf{rw}}(n) \leq n \cdot \mathsf{oc}(\alpha^*(n))$.

Putting these together with a union bound, the result follows.

The second part of the theorem follows from substituting in the specific choice of $\gamma$. $\qquad\square$

# 5  Applying to STORM and SASS

In this section, we demonstrate how the generic oracle complexity bounds in the previous section can be applied to concrete combinations of oracles and algorithms. We will consider the specific setting of expected risk minimization and two algorithms, first-order STORM and SASS, which are described earlier in the paper and fully analyzed in [BCMS19] and [JSX21a], respectively. For each case, we will state the bounds on $\mathsf{oc}(\alpha)$ as a function of $\alpha$, and use those bounds in conjunction with the known bounds on $T_\varepsilon$ (that have been derived in previous papers), to obtain a bound on the total oracle complexity for each algorithm.

The results we obtain are the first ones that bound the total oracle complexity of STORM and SASS, and we show that both algorithms are essentially near optimal in terms of total gradient sample complexity. When deriving these results, for simplicity of presentation, we omit most of the constants involved in the specific bounds on $T_\varepsilon$ and specific conditions on various algorithmic constants. For all such details, we refer the reader to [BCMS19] and [JSX21a]. We will include short comments regarding these constants, but otherwise replace them with a $\mathcal{O}(\cdot)$ notation.

**Problem Setting: Expected risk minimization (ERM)** can be written as

$$\min_{x \in \mathbb{R}^m} \phi(x) = \mathbb{E}_{d \sim \mathcal{D}}[\ell(x, d)].$$

Here, $x$ represents the vector of model parameters, $d$ is a data sample following distribution $\mathcal{D}$, and $\ell(x, d)$ is the loss when the model parameterized by $x$ is evaluated on data point $d$. This problem is central in supervised machine learning and other settings such as simulation optimization [KPH15]. For this problem, it is common to assume the function $\phi$ is $L$-smooth and is bounded from below, and gradients of functions $\nabla_x \ell(x, d)$ can be computed for any $d \sim \mathcal{D}$, so we will consider this setting in this section.

In this setting, the zeroth- and first-order oracles are usually computed as follows:

$$f(x, \mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{d \in \mathcal{S}} \ell(x, d), \quad g(x, \mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{d \in \mathcal{S}} \nabla_x \ell(x, d), \tag{9}$$

where $\mathcal{S}$ is the "minibatch" - that is a set of i.i.d samples from $\mathcal{D}$. Generally, $|\mathcal{S}|$ can be chosen to depend on $x$.

In what follows we will refer to the total number of times an algorithm computes $\ell(x, d)$ for a specific $x$ and $d$ as its *total function value sample complexity* and the number of times the algorithm computes $\nabla \ell(x, d)$ as its *total gradient sample complexity*. The total (oracle or sample) complexity of the algorithm is defined as the sum of these two quantities.

## 5.1   Total sample complexity of first-order STORM

We first consider the first-order stochastic trust-region method (STORM) as introduced and analyzed in [BCMS19]. The algorithm uses zeroth- and first-order oracles defined in TR1.0 and TR1.1 in Section 2. Trust-region algorithms are usually applied to nonconvex functions and the stopping time of STORM is defined as $T_\varepsilon = \min\{k : \|\nabla \phi(x_k)\| \leq \varepsilon\}$. In Section 3.3 of [BCMS19], it is shown that Assumption 1 is satisfied with $\bar{\alpha} = \frac{\varepsilon}{\zeta}$, where $\zeta$ is a moderate constant that depends on $\kappa_{eg}$, $L$ and some constant chosen by the algorithm.

In [BCMS19], the oracle costs of STORM in the ERM setting are briefly discussed under the following **assumptions** on $\ell(x, d)$.

- Function value: It is assumed that there is some $\sigma_f \geq 0$ such that for all $x$, $\mathrm{Var}_{d \sim \mathcal{D}} [\ell(x, d)] \leq \sigma_f^2$.

- Gradient: It is assumed that $\mathbb{E}_{d \sim \mathcal{D}}[\nabla_x \ell(x, d)] = \nabla \phi(x)$, and that there is some $\sigma_g \geq 0$ such that for all $x$,

$$\mathbb{E}_{d \sim \mathcal{D}} \|\nabla_x \ell(x, d) - \nabla \phi(x)\|^2 \leq \sigma_g^2. \tag{10}$$

The cost of each oracle call is the number of samples in the associated minibatch $\mathcal{S}$. By applying Chebyshev's inequality it is easy to bound the oracle costs of TR1.0 and TR1.1.

- Cost of TR1.0 with parameter $\alpha$: $\mathsf{oc}_0(\alpha) = \frac{\sigma_f^2}{\delta_0 \kappa_{ef}^2 \alpha^4}$,

- Cost of TR1.1 with parameter $\alpha$: $\mathsf{oc}_1(\alpha) = \frac{\sigma_g^2}{\delta_1 \kappa_{eg}^2 \alpha^2}$.

Below we substitute the specific oracle costs into Theorem 2 to obtain the expected total sample complexity for the first-order STORM algorithm. Specifically, we will bound the total sample complexity of STORM $\mathbb{E}[\mathsf{TOC}(\min\{T_\varepsilon, n\})]$ by deriving bounds on the expectation of the total function value sample complexity $\mathsf{TOC}_0$ and the total gradient sample complexity $\mathsf{TOC}_1$, where

$$\mathsf{TOC}_0(n) = \sum_{k=1}^{n} \mathsf{oc}_0(\mathcal{A}_k), \ \ \mathsf{TOC}_1(n) = \sum_{k=1}^{n} \mathsf{oc}_1(\mathcal{A}_k) \ \ \text{and} \ \ \mathsf{TOC}(n) = \mathsf{TOC}_0(n) + \mathsf{TOC}_1(n).$$

**Theorem 4** (Expected Total Sample Complexity Bound of First-Order STORM). *Let $p = 1 - \delta_0 - \delta_1$ and $q = 1 - p$. For the first-order STORM algorithm, for any iteration number $n \in \mathbb{Z}^+$, and $\gamma > (2q)^{\frac{1}{4}}$, we have*

$$\mathbb{E}[\mathsf{TOC}(\min\{T_\varepsilon, n\})] \leq \mathcal{O}\left(n \log_{p/q}(n) \left(\frac{\sigma_f^2}{\varepsilon^4} n^{4 \log_{q/p} \gamma} + \frac{\sigma_g^2}{\varepsilon^2} n^{2 \log_{q/p} \gamma}\right)\right).$$

*If $\gamma \geq \left(\frac{q}{p}\right)^{\frac{\log c}{\log n}}$ for some constant $c > 1$ (so that $n^{\log_{q/p} \gamma} \leq c$), the above simplifies to be*

$$\mathbb{E}[\mathsf{TOC}(\min\{T_\varepsilon, n\})] \leq n \log(n) \cdot \mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^4} + \frac{\sigma_g^2}{\varepsilon^2}\right). \tag{11}$$

*Proof.* By Theorem 2, the total expected cost of the zeroth-order oracle over $n$ iterations is bounded above by:

$$\mathbb{E}[\mathsf{TOC}_0(\min\{T_\varepsilon, n\})] \leq n \sum_{\ell=1}^{n} \min\left\{1, n\left(\frac{q}{p}\right)^\ell + \frac{2\sqrt{pq}}{(1-2\sqrt{pq})^2}(2q)^\ell\right\} \cdot \mathsf{oc}_0(\bar{\alpha}\gamma^\ell) + n\,\mathsf{oc}_0(\bar{\alpha})$$

$$\leq n \underbrace{\sum_{\ell=1}^{n} \min\left\{1, n\left(\frac{q}{p}\right)^\ell\right\} \cdot \mathsf{oc}_0(\bar{\alpha}\gamma^\ell)}_{=:A}$$

$$+ \frac{2n\sqrt{pq}}{(1-2\sqrt{pq})^2} \underbrace{\sum_{\ell=1}^{n}(2q)^\ell \cdot \mathsf{oc}_0(\bar{\alpha}\gamma^\ell)}_{=:B} + n\,\mathsf{oc}_0(\bar{\alpha}).$$

For the zeroth-order oracle, $\mathsf{oc}_0(\alpha) = \frac{\sigma_f^2}{\delta_0 \kappa_{ef}^2 \alpha^4} = \mathcal{O}(\frac{\sigma_f^2}{\alpha^4})$. We use this to calculate upper bounds for $A$ and $B$. First, we consider $A$. Note that $\min\{1, n(\frac{q}{p})^\ell\} = 1$, if and only if $\ell \leq \log_{p/q}(n)$. Therefore,

$$A \leq \sum_{\ell=1}^{\log_{\frac{p}{q}}(n)} \mathsf{oc}_0(\bar{\alpha}\gamma^\ell) + n \sum_{\ell \geq \log_{\frac{p}{q}}(n)} \left(\frac{q}{p}\right)^\ell \mathsf{oc}_0(\bar{\alpha}\gamma^\ell)$$

$$\leq \sum_{\ell=1}^{\log_{\frac{p}{q}}(n)} \mathcal{O}\left(\frac{\sigma_f^2}{\bar{\alpha}^4\gamma^{4\ell}}\right) + n \sum_{\ell \geq \log_{\frac{p}{q}}(n)} \left(\frac{q}{p}\right)^\ell \mathcal{O}\left(\frac{\sigma_f^2}{\bar{\alpha}^4\gamma^{4\ell}}\right)$$

$$\leq \mathcal{O}\left(\frac{\sigma_f^2}{\bar{\alpha}^4}\right) \left(\sum_{\ell=1}^{\log_{\frac{p}{q}}(n)} \frac{1}{\gamma^{4\ell}} + n \sum_{\ell \geq \log_{\frac{p}{q}}(n)} \left(\frac{q}{p}\right)^\ell \frac{1}{\gamma^{4\ell}}\right)$$

$$\leq \mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^4}\right) \left(\log_{\frac{p}{q}}(n) \cdot \left(\frac{1}{\gamma^4}\right)^{\log_{\frac{p}{q}}(n)} + n \left(\frac{q}{p\gamma^4}\right)^{\log_{\frac{p}{q}}(n)} \frac{1}{1 - \frac{q}{p\gamma^4}}\right)$$

$$= \mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^4} \log_{\frac{p}{q}}(n)\, n^{\log_{\frac{p}{q}}\left(\frac{1}{\gamma^4}\right)}\right)$$

$$= \mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^4} \log_{\frac{p}{q}}(n)\, n^{4\log_{\frac{q}{p}}\gamma}\right).$$

Next we bound $B$. We have

$$B = \sum_{\ell=1}^{n}(2q)^\ell \cdot \mathsf{oc}_0(\bar{\alpha}\gamma^\ell) \leq \sum_{\ell=0}^{\infty}(2q)^\ell\, \mathcal{O}\left(\frac{\sigma_f^2}{\bar{\alpha}^4\gamma^{4\ell}}\right) \leq \mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^4}\right)\left(\frac{1}{1 - 2q \cdot \frac{1}{\gamma^4}}\right) = \mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^4}\right).$$

Using these bounds on $A$ and $B$ in the expression for $\mathbb{E}[\mathsf{TOC}_0(\min\{T_\epsilon, n\})]$, we obtain the bound on the total function value sample complexity as $\mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^4} n \log_{\frac{p}{q}}(n)\, n^{4\log_{\frac{q}{p}}\gamma}\right)$.

A similar calculation using the cost of the first-order oracle yields the bound $\mathcal{O}\left(\frac{\sigma_g^2}{\varepsilon^2} n \log_{\frac{p}{q}}(n)\, n^{2\log_{\frac{q}{p}}\gamma}\right)$ for $\mathbb{E}[\mathsf{TOC}_1(\min\{T_\epsilon, n\})]$. Since $\mathsf{TOC}(\min\{T_\epsilon, n\}) = \mathsf{TOC}_0(\min\{T_\epsilon, n\}) + \mathsf{TOC}_1(\min\{T_\epsilon, n\})$ by definition, the result follows.

$\square$

Let us discuss the implications of Theorem 4. In [ACD+19], a lower bound on the total gradient sample complexity for stochastic optimization of non-convex, smooth functions is derived and shown to be, in the worst case, $\frac{C}{\varepsilon^4}$, for some positive constant $C$. This complexity lower bound holds even when exact function

values $\phi(x)$ are also available. We note that the definition of complexity in [ACD$^+$19] is the smallest number of sample gradient evaluations required to return a point $x$ with $\mathbb{E}[\|\nabla\phi(x)\|] \leq \varepsilon$, which is different from $\mathsf{TOC}(T_\epsilon)$ which we are aiming to bound here. We believe that the lower bound in [ACD$^+$19] applies to our definition as well, but this is a subject of a separate study.

In [BCMS19], it is shown that $\mathbb{E}[T_\varepsilon] \leq \frac{C_1}{\varepsilon^2}$ for some $C_1$ sufficiently large that depends on $\delta_1, \delta_0, \kappa_{eg}, L$ and some algorithmic constants. Thus, if $n = \frac{C_1}{\varepsilon^2}$ in inequality (11) of Theorem 4 , as long as $\gamma$ is sufficiently large, we obtain $\mathbb{E}[\mathsf{TOC}(\min\{T_\varepsilon, n\})] \leq \mathcal{O}\left(\left(\frac{\sigma_f^2}{\varepsilon^6} + \frac{\sigma_g^2}{\varepsilon^4}\right)\log(\frac{1}{\varepsilon})\right)$. In particular, the total gradient sample complexity is $\mathcal{O}\left(\frac{\sigma_g^2}{\varepsilon^4}\log\left(\frac{1}{\varepsilon}\right)\right)$, which essentially matches the complexity lower bound as described in [ACD$^+$19] up to a logarithmic factor. The total function value sample complexity is worse than that of the gradient if $\sigma_f^2$ is large. However, if $\sigma_f^2 \leq \mathcal{O}(\sigma_g^2\varepsilon^2)$ (which often happens in practice since $\sigma_g$ usually scales with the dimension of the problem, and tends to be much larger than $\sigma_f$), the total sample complexity bound of STORM matches the lower bound up to a logarithmic factor.

We note now that choosing $n = \frac{C_1}{\varepsilon^2}$ in Theorem 4 does not in fact guarantee that $T_\varepsilon < n$, since for STORM, only a bound on $\mathbb{E}[T_\varepsilon]$ has been derived. However, this statement can be made true in probability, thanks to Theorem 3, by simply applying Markov inequality for $n = C_2\frac{C_1}{\varepsilon^2}$ (where $C_2 > 1$).

**Theorem 5** (High Probability Total Sample Complexity Bound of First-Order STORM)**.** *For the first-order STORM algorithm applied to expected risk minimization, let $n$ be chosen such that $n \geq C_2\frac{C_1}{\varepsilon^2}$ (for some $C_1$ sufficiently large so that $\frac{C_1}{\varepsilon^2} \geq \mathbb{E}[T_\varepsilon]$, and any $C_2 > 1$), and $\gamma$ be chosen so that $\gamma \geq \max\left\{\frac{1}{2}, \left(\frac{1}{2q}\right)^{\frac{\log(2\beta)}{(1+\omega)\log n}}\right\}$ (for some $\beta \leq \frac{1}{2}$, and any $\omega > 0$). Then, with probability at least $1 - \frac{1}{C_2} - \mathcal{O}(n^{-\omega})$,*

$$\mathsf{TOC}(T_\varepsilon) \leq \mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^6} + \frac{\sigma_g^2}{\varepsilon^4}\right). \tag{12}$$

*Proof.* The theorem is a simple application of Theorem 3 to the specific setting. $\square$

**Remark 2.** *1. Compared to the expected total sample complexity bound, this high probability bound is smaller by a log factor.*

*2. In [GRVZ18], a first-order trust-region algorithm similar to STORM with the same first-order oracle (i.e. TR1.1 with $\epsilon_g = 0$), but with an exact zeroth-order oracle (i.e. TR1.0 with $\kappa_{ef} = \delta_0 = 0$) is analyzed. In this case, it is shown that $\mathbb{P}(T_\varepsilon > n) \leq \exp(-C_1 n)$ (for some constant $C_1$ depends on $\delta_1$), for any $n \geq \frac{C_2}{\varepsilon^2}$ (with some sufficiently large $C_2$). Using a similar application of Theorem 3, we can show that as long as $\gamma$ is sufficiently large, the total gradient sample complexity of that trust-region algorithm is bounded above by $\mathcal{O}(\frac{\sigma_g^2}{\varepsilon^4})$ with probability at least $1 - \exp(-C_1 n) - \mathcal{O}(n^{-\omega})$ (which is a significant improvement over the probability in Theorem 5).*

*3. Another first-order trust-region algorithm, with weaker oracle assumptions than those in [GRVZ18] is introduced and analyzed in [CBS22]. This algorithm relies on the first-order oracle as described in TR1.1 and the zeroth-order oracle as described in SS.0. For this algorithm, it is shown that $\mathbb{P}(T_\varepsilon > n) \leq 2\exp(-C_1 n) + \exp(-C_2)$ ($C_2$ being any positive constant), where $n = C_3\frac{C_2}{\varepsilon^2}$ for some sufficiently large $C_3$ and some positive $C_1$ that depends on $\delta_0$ and $\delta_1$. Thus, again, using Theorem 3 we can show that as long as $\gamma$ is sufficiently large, the total sample complexity of the first-order trust-region algorithm in [CBS22] is bounded above by $\mathcal{O}(\frac{\sigma_f^2}{\varepsilon^6} + \frac{\sigma_g^2}{\varepsilon^4})$ with probability at least $1 - 2\exp(-C_1 n) - \exp(-C_2) - \mathcal{O}(n^{-\omega})$.*

## 5.2 Total sample complexity of SASS

We now consider the SASS algorithm[1], analyzed in [BCS19, JSX21a] and described in Section 2.3. By Proposition 1, 2 and 4 of [JSX21a], Assumption 1 is satisfied, with $\bar{\alpha}$ as given in the propositions.

In the empirical risk minimization setting, the following **assumptions** on $\ell(x, d)$ are made in [JSX21a].

---

[1]This algorithm was also referred to as ALOE in [JSX21b]. Its name has been changed to SASS since [JSX21a].

- Function value: It is assumed that $|\ell(x,d) - \phi(x)|$ is a subexponential random variable and that there is some $\sigma_f \geq 0$ such that for all $x$, $\mathrm{Var}_{d \sim \mathcal{D}} [\ell(x,d)] \leq \sigma_f^2$.
  For example, if $\ell(x,d)$ is uniformly bounded, then $|\ell(x,d) - \phi(x)|$ is subexponential.

- Gradient: It is assumed that $\mathbb{E}_{d \sim \mathcal{D}}[\nabla_x \ell(x,d)] = \nabla\phi(x)$, and that for some $M_c, M_v \geq 0$ and for all $x$,

$$\mathbb{E}_{d \sim \mathcal{D}} \|\nabla_x \ell(x,d) - \nabla\phi(x)\|^2 \leq M_c + M_v \|\nabla\phi(x)\|^2. \tag{13}$$

This assumption is fairly general and is studied in the literature [BCN18].

For non-convex functions, the stopping time is defined as $T_\varepsilon = \min\{k : \|\nabla\phi(x_k)\| \leq \varepsilon\}$, same as in the case of STORM. For strongly convex functions, the stopping time is defined as $T_\varepsilon = \min\{k : \phi(x_k) - \inf_x \phi(x) \leq \varepsilon\}$. To achieve the desired accuracy, oracles SS.0 and SS.1 have to be sufficiently accurate in the sense that $\epsilon_f$ and $\epsilon_g$ have to be sufficiently small with respect to $\varepsilon$. In the case of expected risk minimization, the oracles can be implemented for any $\epsilon_f$ and $\epsilon_g$ by choosing an appropriate mini-batch size. Thus, here we will first fix $\varepsilon$ and then discuss the oracles that deliver sufficient accuracy for such $\varepsilon$, for the theory in [JSX21a] to apply.

**Oracle Costs per Iteration.** In [JSX21a], it is shown that given the desired convergence tolerance $\varepsilon$, sufficiently accurate oracles SS.0 and SS.1 can be implemented for any step size parameter $\alpha$ as follows:

- Zeroth-order oracle: Proposition 5 of [JSX21a] shows that a sufficiently accurate zeroth-order oracle can be obtained by using a minibatch of size

$$\mathsf{oc}_0(\alpha) = \begin{cases} \mathcal{O}(\sigma_f^2/\varepsilon^4), & \text{for the non-convex case,} \\ \mathcal{O}(\sigma_f^2/\varepsilon^2), & \text{for the strongly convex case.} \end{cases} \tag{14}$$

Note that the cost of the zeroth-order oracle is independent of $\alpha$.

- First-order oracle: Proposition 6 of [JSX21a] implies that a sufficiently accurate first-order oracle can be obtained by using a minibatch of size

$$\mathsf{oc}_1(\alpha) = \begin{cases} \mathcal{O}\left(\frac{M_c}{\varepsilon^2} + \frac{M_v}{\min\{\tau, \kappa\alpha\}^2}\right), & \text{for the non-convex case,} \\ \mathcal{O}\left(\frac{M_c}{\varepsilon} + \frac{M_v}{\min\{\tau, \kappa\alpha\}^2}\right), & \text{for the strongly convex case.} \end{cases} \tag{15}$$

The cost of the first-order oracle is indeed non-increasing in $\alpha$, so Assumption 2 is satisfied. For simplicity of the presentation and essentially without loss of generality, we will assume $\tau \geq \kappa\bar{\alpha}$.

Substituting these bounds into Theorem 2, we obtain the following expected total sample complexity.

**Theorem 6** (Expected Total Sample Complexity of SASS)**.** *For the SASS algorithm applied to expected risk minimization, for any iteration number $n \in \mathbb{Z}^+$, and any $\gamma > (2q)^{\frac{1}{2}}$, we have*

- *Non-convex case:*

$$\mathbb{E}[\mathsf{TOC}(\min\{T_\varepsilon, n\})] \leq \mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^4} \cdot n + \frac{M_c}{\varepsilon^2} \cdot n + M_v \cdot n \left(n^{\log_{\frac{p}{q}}\left(\frac{1}{\gamma^2}\right)} \log_{\frac{p}{q}}(n)\right)\right).$$

- *Strongly convex case:*

$$\mathbb{E}[\mathsf{TOC}(\min\{T_\varepsilon, n\})] \leq \mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^2} \cdot n + \frac{M_c}{\varepsilon} \cdot n + M_v \cdot n \left(n^{\log_{\frac{p}{q}}\left(\frac{1}{\gamma^2}\right)} \log_{\frac{p}{q}}(n)\right)\right).$$

*Moreover, if $\gamma \geq \left(\frac{q}{p}\right)^{\frac{\log c}{2\log n}}$ for some constant $c > 1$ (so that $n^{\log_{\frac{p}{q}}\left(\frac{1}{\gamma^2}\right)} \leq c$), the above simplifies to*

- *Non-convex case:*

$$\mathbb{E}[\mathsf{TOC}(\min\{T_\epsilon, n\})] \leq \mathcal{O}\left(n\left(\frac{\sigma_f^2}{\varepsilon^4} + \frac{M_c}{\varepsilon^2} + M_v \log(n)\right)\right). \tag{16}$$

- *Strongly convex case:*

$$\mathbb{E}[\mathsf{TOC}(\min\{T_\epsilon, n\})] \leq \mathcal{O}\left(n\left(\frac{\sigma_f^2}{\varepsilon^2} + \frac{M_c}{\varepsilon} + M_v \log(n)\right)\right). \tag{17}$$

*Proof.* Since the cost of each call to the zeroth-order oracle (14) is independent of $\alpha$, the total function value sample complexity over $n$ iterations is simply obtained by multiplying (14) by $n$.

The cost of the first-order oracle (15) consists of two parts, $\mathsf{oc}_1(\alpha) = \mathsf{oc}_{1,1}(\alpha) + \mathsf{oc}_{1,2}(\alpha)$. The first part, $\mathsf{oc}_{1,1}(\alpha)$, is $\mathcal{O}(\frac{M_c}{\varepsilon^2})$. Since this is *independent* of $\alpha$, the total contribution of this part to the total gradient sample complexity over $n$ iterations is $n\,\mathsf{oc}_{1,1}(\alpha)$, which is bounded by $\mathcal{O}(\frac{M_c}{\varepsilon^2}\,n)$.

The second part of the cost of the first-order oracle is $\mathsf{oc}_{1,2}(\alpha) := \mathcal{O}(\frac{M_v}{\min\{\tau, \kappa\alpha\}^2})$. By Theorem 2, the total expected cost over $n$ iterations from this part is bounded above by:

$$n\sum_{\ell=1}^{n}\min\left\{1, n\left(\frac{q}{p}\right)^\ell + \frac{2\sqrt{pq}}{(1-2\sqrt{pq})^2}(2q)^\ell\right\} \cdot \mathsf{oc}_{1,2}(\bar\alpha\gamma^\ell) + n\,\mathsf{oc}_{1,2}(\bar\alpha)$$

$$\leq n\underbrace{\sum_{\ell=1}^{n}\min\left\{1, n\left(\frac{q}{p}\right)^\ell\right\} \cdot \mathsf{oc}_{1,2}(\bar\alpha\gamma^\ell)}_{=:A} + \frac{2n\sqrt{pq}}{(1-2\sqrt{pq})^2}\underbrace{\sum_{\ell=1}^{n}(2q)^\ell \cdot \mathsf{oc}_{1,2}(\bar\alpha\gamma^\ell)}_{=:B} + n\mathsf{oc}_{1,2}(\bar\alpha).$$

Note that the expression above only involves $\mathsf{oc}_{1,2}(\alpha)$ for $\alpha \leq \bar\alpha$. Therefore, by our earlier assumption that $\tau \geq \kappa\bar\alpha$, we have $\mathsf{oc}_{1,2}(\alpha) = \mathcal{O}(\frac{M_v}{\kappa^2\alpha^2}) = \mathcal{O}(\frac{M_v}{\alpha^2})$ (since $\kappa$ is a constant). We now use this to calculate upper bounds for $A$ and $B$. Using similar arguments as the proof of Theorem 4, we have

$$A = \mathcal{O}\left(\frac{M_v}{\bar\alpha^2}\log_{\frac{p}{q}}(n)\,n^{\log_{\frac{p}{q}}\left(\frac{1}{\gamma^2}\right)}\right) \text{ and } B = \mathcal{O}\left(\frac{M_v}{\bar\alpha^2}\right).$$

Together with the previous arguments, the result follows. $\qquad\square$

In [JSX21a], the iteration bound on $T_\epsilon$ is shown to be $\frac{C_1}{\varepsilon^2} + \log_{1/\gamma}\frac{\alpha_0}{\bar\alpha}$ in the non-convex case, and $C_2\log\frac{1}{\varepsilon} + \log_{1/\gamma}\frac{\alpha_0}{\bar\alpha}$ in the strongly convex case (for some $C_1$ and $C_2$ sufficiently large) in high probability. Thus, we can select $n$ appropriately and derive the high probability bound on the total sample complexity of SASS using Theorem 3.

**Theorem 7** (High Probability Total Sample Complexity of SASS). *For the SASS algorithm applied to expected risk minimization, let $n \geq \frac{C_1}{\varepsilon^2} + \log_{1/\gamma}\frac{\alpha_0}{\bar\alpha}$ in the non-convex case, and $n \geq C_2\log\frac{1}{\varepsilon} + \log_{1/\gamma}\frac{\alpha_0}{\bar\alpha}$ in the strongly convex case.*

*For any $\omega > 0$, with probability at least $1 - 2\exp(-C_3 n) - \mathcal{O}(n^{-\omega})$ (for some $C_3 > 0$), we have*

- *Non-convex case:*

$$\mathsf{TOC}(T_\varepsilon) \leq \mathcal{O}\left(\left(\frac{1}{\varepsilon^2} + \log_{1/\gamma}\frac{\alpha_0}{\bar\alpha}\right) \cdot \left(\frac{\sigma_f^2}{\varepsilon^4} + \frac{M_c}{\varepsilon^2} + \frac{M_v}{\varepsilon^{4(1+\omega)\log_{2q}(\gamma)}}\right)\right). \tag{18}$$

- *Strongly convex case:*

$$\mathsf{TOC}(T_\varepsilon) \leq \mathcal{O}\left(\left(\log\frac{1}{\varepsilon} + \log_{1/\gamma}\frac{\alpha_0}{\bar\alpha}\right) \cdot \left(\frac{\sigma_f^2}{\varepsilon^2} + \frac{M_c}{\varepsilon} + M_v\left(\log\frac{1}{\varepsilon}\right)^{2(1+\omega)\log_{2q}(\gamma)}\right)\right). \tag{19}$$

*If $\gamma \in \left[\max\left\{\frac{1}{2}, \left(\frac{1}{2q}\right)^{\frac{\log(2\beta)}{(1+\omega)\log n}}\right\}, \left(\frac{\bar\alpha}{\alpha_0}\right)^{\frac{1}{cn}}\right]$, where $\beta$ is any constant smaller than $\frac{1}{2}$, and $c$ is any constant in $(0,1)$, the above simplifies to*

- *Non-convex case:*

$$\mathsf{TOC}(T_\varepsilon) \leq \mathcal{O}\left(\frac{1}{\varepsilon^2} \cdot \left(\frac{\sigma_f^2}{\varepsilon^4} + \frac{M_c}{\varepsilon^2} + \frac{M_v}{\beta^2}\right)\right). \tag{20}$$

- *Strongly convex case:*

$$\mathsf{TOC}(T_\varepsilon) \leq \mathcal{O}\left(\log\frac{1}{\varepsilon} \cdot \left(\frac{\sigma_f^2}{\varepsilon^2} + \frac{M_c}{\varepsilon} + \frac{M_v}{\beta^2}\right)\right). \tag{21}$$

*Proof.* The bounds (18) and (19) follow by using (14) and (15) in Theorem 3, and Theorem 3.8 of [JSX21a].

The bounds (20) and (21) follow from (18) and (19), respectively, by using the fact that $\gamma$ lies in the appropriate range and $\log_{1/\gamma}\frac{\alpha_0}{\alpha} \leq cn$ with $c \in (0,1)$.

In the non-convex case, we have $\frac{1}{\varepsilon^2} + \log_{1/\gamma}\frac{\alpha_0}{\alpha} = \mathcal{O}(\frac{1}{\varepsilon^2})$ and $\frac{M_v}{\varepsilon^{4(1+\omega)\log_{2q}(\gamma)}} = \mathcal{O}(\frac{M_v}{\beta^2})$ when $\gamma$ lies in the specified range. It is worth noting that $c \in (0,1)$ implies there is some $n = \mathcal{O}(\frac{1}{\varepsilon^2})$ that satisfies $n \geq \frac{C_1}{\varepsilon^2} + \log_{1/\gamma}\frac{\alpha_0}{\alpha}$. The result for the strongly convex case follows similarly. $\square$

**Remark 3.** *1. We consider some of the implications of Theorem 7 below. Similar implications hold for the expected total sample complexity.*

- *In the non-convex case, from (20), the total function value sample complexity is $\mathcal{O}\left(\frac{\sigma_f^2}{\varepsilon^6}\right)$ and the total gradient sample complexity is $\mathcal{O}\left(\frac{M_c}{\varepsilon^4} + \frac{M_v}{\varepsilon^2}\right)$. In particular, the total gradient sample complexity matches that of SGD, and it essentially matches the complexity lower bound as described in [ACD+19] (for a different definition of complexity). Specifically, if $\sigma_f = 0$ (i.e., function values are exact), the lower bound in [ACD+19] applies and the total sample complexity of SASS matches it.*

- *If $M_c = 0$ (sometimes referred to as the interpolation case), then the total gradient sample complexity reduces to $\mathcal{O}\left(\frac{M_v}{\varepsilon^2}\right)$. Hence, the total gradient sample complexity matches that of SGD under interpolation [BCN18, BM11].*

- *In the strongly convex case, the total function value sample complexity is $\mathcal{O}(\frac{\sigma_f^2}{\varepsilon^2}\log\frac{1}{\varepsilon})$ and the total gradient sample complexity is $\mathcal{O}((\frac{M_c}{\varepsilon} + M_v)\log\frac{1}{\varepsilon})$. In particular, the total gradient sample complexity matches that of SGD up to a logarithmic term.*

*2. Our framework can be applied to the convex setting as well. We focus on the non-convex and strongly convex settings in this paper for brevity and to cleanly illustrate how our framework can be applied, since the convex case requires some additional complications in the presentation. These complications are present in the previous papers that analyze iteration complexity bounds, and are not specific to this paper. For example, the stopping time in the convex setting for SASS is defined in terms of two parameters $\varepsilon = (\varepsilon_0, \varepsilon_1)$ as follows: $T_\varepsilon = \min\{k : \phi(X_k) - \phi(x^*) \leq \varepsilon_0 \text{ or } \|\nabla\phi(X_k)\| \leq \varepsilon_1\}$.*

# 6   Conclusion

We analyzed the behavior of the step size parameter in Algorithm 1, an adaptive stochastic optimization framework that encompasses a wide class of algorithms analyzed in recent literature. We derived a high probability lower bound for this parameter, and as a result, developed a simple strategy for controlling this lower bound.

For many settings, having a fixed lower bound on the step size parameter implies an upper bound on the cost of the oracles that compute the gradient and function estimates. We developed a framework to analyze the expected and high probability total oracle complexity for this general class of algorithms, and illustrated the use of it by deriving total sample complexity bounds for two specific algorithms - the first-order stochastic trust-region (STORM) algorithm [BCMS19] and a stochastic step search (SASS) algorithm [JSX21a] in the expected risk minimization setting. We showed that the sample complexity of both these algorithms essentially matches the complexity lower bound of first-order algorithms for stochastic non-convex optimization [ACD+19], which was not known before.

# Acknowledgments

# References

[ACD+19]   Yossi Arjevani, Yair Carmon, John C Duchi, Dylan J Foster, Nathan Srebro, and Blake Wood-worth. Lower bounds for non-convex stochastic optimization. *arXiv preprint arXiv:1912.02365*, 2019.

[Ans60]   Francis John Anscombe. Rejection of outliers. *Technometrics*, 2:123–146, 1960.

[BCCS21]   Albert S Berahas, Liyuan Cao, Krzysztof Choromanski, and Katya Scheinberg. A theoretical and empirical comparison of gradient approximations in derivative-free optimization. *Foundations of Computational Mathematics*, 2021.

[BCMS19]   Jose Blanchet, Coralia Cartis, Matt Menickelly, and Katya Scheinberg. Convergence rate analysis of a stochastic trust-region method via supermartingales. *INFORMS journal on optimization*, 1(2):92–119, 2019.

[BCN18]   Léon Bottou, Frank E Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *Siam Review*, 60(2):223–311, 2018.

[BCS19]   Albert S Berahas, Liyuan Cao, and Katya Scheinberg. Global convergence rate analysis of a generic line search algorithm with noise. *SIAM Journal on Optimization*, 2019.

[BM11]   F. Bach and E. Moulines. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain.*, pages 451–459, 2011.

[BXZ23]   Albert S. Berahas, Miaolan Xie, and Baoyu Zhou. A sequential quadratic programming method with high probability complexity bounds for nonlinear equality constrained stochastic optimization, 2023.

[CBS22]   Liyuan Cao, Albert S Berahas, and Katya Scheinberg. First-and second-order high probability complexity bounds for trust-region methods with noisy oracles. *arXiv preprint arXiv:2205.03667*, 2022.

[CS17]   C. Cartis and K. Scheinberg. Global convergence rate analysis of unconstrained optimization methods based on probabilistic models. *Mathematical Programming*, 169(2):337–375, 2017.

[Fel68]   William Feller. *An Introduction to Probability Theory and Its Applications*, volume 1. Wiley, January 1968.

[GRVZ18]   Serge Gratton, Clément W Royer, Luís N Vicente, and Zaikun Zhang. Complexity and global rates of trust-region methods based on probabilistic models. *IMA Journal of Numerical Analysis*, 38(3):1579–1597, 2018.

[JSX21a]   Billy Jin, Katya Scheinberg, and Miaolan Xie. High probability complexity bounds for adaptive step search based on stochastic oracles. *arXiv preprint arXiv:2106.06454*, 2021.

[JSX21b]   Billy Jin, Katya Scheinberg, and Miaolan Xie. High probability complexity bounds for line search based on stochastic oracles. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 9193–9203, 2021.

[KB14]     Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[KPH15]    Sujin Kim, Raghu Pasupathy, and Shane G Henderson. A guide to sample average approximation. *Handbook of simulation optimization*, pages 207–243, 2015.

[MWX23]    Matt Menickelly, Stefan M. Wild, and Miaolan Xie. A stochastic quasi-newton method in the absence of common random numbers, 2023.

[PS20]     Courtney Paquette and Katya Scheinberg. A stochastic line search method with expected complexity analysis. *SIAM Journal on Optimization*, 30(1):349–376, 2020.

[RJM17]    Jeffrey Regier, Michael I Jordan, and Jon McAuliffe. Fast black-box variational inference through stochastic trust-region optimization. *Advances in Neural Information Processing Systems*, 30, 2017.

[RVZ23]    F Rinaldi, LN Vicente, and D Zeffiro. Stochastic trust-region and direct-search methods: A weak tail bound condition and reduced sample sizing. 2023.

[SHP18]    Sara Shashaani, Fatemeh S Hashemi, and Raghu Pasupathy. Astro-df: A class of adaptive sampling trust-region algorithms for derivative-free stochastic optimization. *SIAM Journal on Optimization*, 28(4):3145–3176, 2018.

[Spa99]    James C. Spall. Stochastic optimization and the simultaneous perturbation method. In *Proceedings of the 31st Conference on Winter Simulation: Simulation—a Bridge to the Future - Volume 1*, WSC '99, page 101–109, New York, NY, USA, 1999. Association for Computing Machinery.

[SX23]     Katya Scheinberg and Miaolan Xie. Stochastic adaptive regularization method with cubics: A high probability complexity bound. In *2023 Winter Simulation Conference (WSC)*, To appear, 2023.

[TMDQ16]   Conghui Tan, Shiqian Ma, Yu-Hong Dai, and Yuqiu Qian. Barzilai-borwein step size for stochastic gradient descent. *Advances in neural information processing systems*, 29, 2016.

[YCKB18]   Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In *International Conference on Machine Learning*, pages 5650–5659. PMLR, 2018.