

RESEARCH OUTPUTS / RÉSULTATS DE RECHERCHE

Fast Stochastic Second-Order Adagrad for Nonconvex Bound-Constrained Optimization

Bellavia, Stefania; Gratton, Serge; Morini, Benedetta; Toint, Philippe

Publication date:
2025

[Link to publication](#)

Citation for published version (HARVARD):

Bellavia, S, Gratton, S, Morini, B & Toint, P 2025 'Fast Stochastic Second-Order Adagrad for Nonconvex Bound-Constrained Optimization' Arxiv.

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Fast Stochastic Second-Order Adagrad for Nonconvex Bound-Constrained Optimization

Stefania Bellavia*, Serge Gratton†, Benedetta Morini‡, Philippe L. Toint§

6 V 2025

Abstract

ADAGB2, a generalization of the **Adagrad** algorithm for stochastic optimization is introduced, which is also applicable to bound-constrained problems and capable of using second-order information when available. It is shown that, given $\delta \in (0, 1)$ and $\epsilon \in (0, 1]$, the ADAGB2 algorithm needs at most $\mathcal{O}(\epsilon^{-2})$ iterations to ensure an ϵ -approximate first-order critical point of the bound-constrained problem with probability at least $1 - \delta$, provided the average root mean square error of the gradient oracle is sufficiently small. Should this condition fail, it is also shown that the optimality level of iterates is bounded above by this average. The relation between the approximate and true classical projected-gradient-based optimality measures for bound constrained problems is also investigated, and it is shown that merely assuming unbiased gradient oracles may be insufficient to ensure convergence in $\mathcal{O}(\epsilon^{-2})$ iterations.

Keywords: Adagrad, stochastic nonlinear optimization, objective-function-free optimization (OFFO), complexity, second-order information, stochastic projected gradients, bound constraints.

1 Introduction

First-order optimization algorithms have been widely used in the contexts on online learning and deep neural network training and their convergence properties on nonconvex problems have been investigated by several authors (see [6] for a survey). Among them, **Adagrad** [16], despite not being always as efficient as others in practice on nonconvex problems [38], enjoys a special position from the theoretical point of view because of its solid and extensive convergence analysis.

When the objective function’s gradient is contaminated by noise (for instance caused by sampling) a probabilistic point of view on the algorithm’s convergence theory is desirable. This has been investigated for **Adagrad** applied to nonconvex functions by a number of authors, as shown in Table 1, along with the relevant important assumptions and results. We discuss this rich body of theory (and the content of this table) in more detail in Section 1.1.

These proposals are however limited, from the theoretician’s point of view, in three respects.

- The first is that the probabilistic analysis is restricted to the case where second-order information is ignored when available. The algorithm remains strictly first-order, with a step always aligned (possibly component-wise) with the approximate steepest-descent direction (no trust region is used). By contrast, other first-order methods have been “augmented” to

*Dipartimento di Ingegneria Industriale, Università degli Studi di Firenze, Firenze, Italia. Member of the INdAM Research Group GNCS. Email: stefania.bellavia@unifi.it.

†Université de Toulouse, INP, IRIT, Toulouse, France. Work partially supported by the Artificial and Natural Intelligence Toulouse Institute (ANITI), French “Investing for the Future - PIA3” program under the Grant agreement ANR-19-PI3A-0004. Email: serge.gratton@enseeiht.fr.

‡Dipartimento di Ingegneria Industriale, Università degli Studi di Firenze, Firenze, Italia. Member of the INdAM Research Group GNCS. Email: benedetta.morini@unifi.it.

§Namur Center for Complex Systems (naXys), University of Namur, Namur, Belgium. Work partially supported by the Artificial and Natural Intelligence Toulouse Institute (ANITI). Email: philippe.toint@unamur.be.

| Paper | | Smooth | Gradients bound / bias | Noise type | 2nd-order usage | Bound constr. | Conv. type | Conv. rate |
|-----------------|------|----------------|---------------------------|-------------------------------------|--------------------|------------------|-----------------------------------|-----------------------------------|
| Li-Orabona | [28] | L^* | no / no | Sub-Gauss. + $\sigma \searrow 0$ | no | no | $\mathbb{E}[\cdot]/\text{w.h.p.}$ | $\mathcal{O}_\ell(\epsilon^{-2})$ |
| Ward et al. | [37] | L | yes / no | Bounded | no | no | w.h.p. | $\mathcal{O}_\ell(\epsilon^{-4})$ |
| Défossez et al. | [15] | L | yes / no | Unrestricted | no | no | $\mathbb{E}[\cdot]$ | $\mathcal{O}_\ell(\epsilon^{-4})$ |
| Gratton et al. | [19] | L | yes / no | Unrestricted | no | no | $\mathbb{E}[\cdot]$ | $\mathcal{O}_\ell(\epsilon^{-4})$ |
| Kavis et al. | [24] | L | yes / no | Sub-Gauss. + $\sigma \searrow 0$ | no | no | w.h.p. | $\mathcal{O}_\ell(\epsilon^{-2})$ |
| Faw et al. | [18] | L | no / no | Affine | no | no | w.h.p. | $\mathcal{O}_\ell(\epsilon^{-4})$ |
| Wang et al. | [34] | (L_0, L_1^*) | no / no | Affine + $\sigma \searrow 0$ | no | no | $\mathbb{E}[\cdot]/\text{w.h.p.}$ | $\mathcal{O}_\ell(\epsilon^{-2})$ |
| Liu et al. | [29] | L | no / no | Sub-Gauss. + $\sigma \searrow 0$ | no | no | w.h.p. | $\mathcal{O}(\epsilon^{-2})$ |
| Attia-Koren | [1] | L | no / no | Affine- + $\sigma \searrow 0$ | no | no | w.h.p. | $\mathcal{O}_\ell(\epsilon^{-2})$ |
| Faw et al. | [17] | (L_0, L_1^*) | no / no | Affine | no | no | w.h.p. | $\mathcal{O}_\ell(\epsilon^{-4})$ |
| Hong-Lin | [22] | L | no / no | Affine- + $\sigma \searrow 0$ | no | no | w.h.p. | $\mathcal{O}_\ell(\epsilon^{-2})$ |
| Hong-Lin | [22] | (L_0, L_1^*) | no / no | Affine-* | no | no | w.h.p. | $\mathcal{O}_\ell(\epsilon^{-4})$ |
| Jiang et al. | [23] | L | no / no | Bounded | no | no | $\mathbb{E}[\cdot]$ | $\mathcal{O}_\ell(\epsilon^{-2})$ |
| This paper | | L | no / no | new1 | yes | yes | $\mathbb{E}[\cdot]/\text{w.h.p.}$ | $\mathcal{O}(\epsilon^{-2})$ |
| This paper | | L | no / yes | new1+new2 | yes | yes | $\mathbb{E}[\cdot]/\text{w.h.p.}$ | $\mathcal{O}(\epsilon^{-2})$ |

Table 1: Convergence theories for the stochastic Adagrad algorithm and their characteristics in the nonconvex setting.

(In the third column, L^* means that the Lipschitz constant must be known, the last column reports a bound (in order) for the considered algorithm to achieve ϵ -criticality, the subscript ℓ indicating the presence of an additional logarithmic factor in ϵ , see Section 1.1 for more detail on the other columns)

use step-sizes along the first-order direction taking second-order information into account (see [32, 2, 4, 13, 14, 35] and the references therein).

- The second is that, at variance with the deterministic case [20], only unconstrained nonconvex problems are considered when the gradient is noisy. While this is adequate for a broad category of applications, it is difficult to use them in a context where *a priori* information on the problem at hand is available, often in the form of constraints, of which bounds on the variables are the most common. Such problems arise, for instance in fracture mechanics [25, 26], inverse problems and PDE-constrained optimization under uncertainty [9]. Moreover, bounds are frequently applied to machine learning models to avoid overfitting [27] or to reflect real-world limits and maintain physical interpretability [30]. Admittedly, these can be taken care of using penalty terms in the loss/objective function, as is often done for Physically Informed Neural Networks (PINNs) (see [7, 36] for instance), but this introduces new hyper-parameters needing calibration and does not ensure that constraints are strictly satisfied. Moreover, a strong penalization of the constraints degrades the problem's conditioning, possibly causing slow convergence, especially if first-order methods are used. We note that [12] have proved convergence (but not complexity) of a stochastic gradient-based interior-point method for bound-constrained problems, but this technique differs considerably from the Adagrad-like gradient algorithms considered here and only allows for diagonal approximations of the Hessian.
- The third shortcoming is that, to the best of the authors' knowledge, bias in the gradient oracle has never been considered for Adagrad, in contrast with other stochastic optimization techniques (see [5, 3] for instance).

The challenge addressed in this paper is thus threefold. We first consider a probabilistic theory of how second-order steps can be accommodated in a stochastic Adagrad-like¹ optimization method.

¹That is an algorithm that reduces to Adagrad when second-order information is not used.

We also show how this can be done in the context of bound-constrained problems, and, finally, propose a theory which does not assume that gradient oracles are unbiased.

More specifically, we consider the problem

$$\min_{x \in \mathcal{F}} f(x) \quad \text{where } \mathcal{F} = \{x \in \mathbb{R}^n \mid l_i \leq x_i \leq u_i\} \quad (1)$$

where f is a smooth (possibly nonconvex) function from an open set containing the feasible region $\mathcal{F} \subseteq \mathbb{R}^n$ into \mathbb{R} , and where l and u are vectors specifying the lower and upper bounds on the variables, respectively (infinite entries in l and/or u are allowed). We also assume that $\nabla_x^1 f(x)$ cannot be computed exactly but is approximated by a random oracle. As a consequence, the algorithm we are about to describe generates a random process, where, for a given iterate x_k , the oracle computes the gradient oracle $g(x_k, \xi)$ where ξ is a random variable (whose distribution may depend on x_k), with probability space $(\Omega, \mathcal{T}_\Omega, \mathbb{P})$. For brevity, we will denote $g_k \stackrel{\text{def}}{=} g(x_k, \xi)$ and we also define $G_k = \nabla_x^1 f(x_k)$ and $H_k = \nabla_x^2 f(x_k)$. The expectation conditioned to knowing g_0, \dots, g_{k-1} will be denoted by the symbol $\mathbb{E}_k[\cdot]$. The symbol $\|\cdot\|$ denotes the 2-norm.

1.1 Related works

The convergence of **Adagrad** in the nonconvex setting has been studied by a number of authors assuming the use of an unbiased stochastic oracle of the true gradient with a variety of assumptions on its noise. In [28], the almost sure asymptotic convergence of the gradients to zero was proved for the first time assuming that the noise has bounded support, i.e.

$$\mathbb{E}_k[\|g_k - G_k\|] \leq \kappa_{\text{bounded}} \quad \kappa_{\text{bounded}} > 0,$$

at any iteration k . Assuming instead that, at each iteration k , the sub-Gaussian noise condition

$$\mathbb{E}_k \left[e^{\|g_k - G_k\|^2 / \sigma^2} \right] \leq e,$$

holds, Li and Orabona also showed that **Adagrad** is noise adaptive, in the sense that, given $\epsilon > 0$, the iteration complexity interpolates, with high probability (denoted by w.h.p. in Table 1), between $\mathcal{O}_\ell(\epsilon^{-4})$ to $\mathcal{O}_\ell(\epsilon^{-2})$ depending on the magnitude of σ , where $\mathcal{O}_\ell(\cdot)$ means that the order is up to a logarithmic term. Unfortunately, this analysis requires *a priori* knowledge of the Lipschitz constant of the gradient in setting the step-size. A similar result was proved in [24] but without assuming this knowledge. Ward, Wu and Bottou in [37] analyzed the convergence of **Adagrad-Norm** assuming uniformly bounded gradients and bounded variance of the gradient oracle at each iteration k , that is

$$\mathbb{E}_k[\|g_k - G_k\|^2] \leq \sigma^2.$$

They showed that the average of the squared norm of the gradients produced by **Adagrad-Norm** converges with high probability at the rate $\mathcal{O}(1/\sqrt{k})$, which implies that the algorithm needs at most $\mathcal{O}_\ell(\epsilon^{-4})$ iterations to achieve $\|G_k\| \leq \epsilon$. Under weaker assumptions, i.e. without assuming boundedness of the variance of the gradient's oracles, these results have been extended (now in expectation form, noted by $\mathbb{E}[\cdot]$ in Table 1) to the component-wise version of **Adagrad** in [15] and to an extended class of methods comprising **Adagrad** in [19].

A different stream of research analyzed **Adagrad** as an optimally-tuned non-adaptive stochastic gradient descent without assuming bounded gradients but under the affine bound on the variance at each iteration given by

$$\mathbb{E}_k[\|g_k - G_k\|^2] \leq \kappa_{1,\text{affine}} + \kappa_{2,\text{affine}} \|G_k\|^2, \quad (2)$$

with $\kappa_{1,\text{affine}} \geq 0$ and $\kappa_{2,\text{affine}} \geq 0$. Under these assumptions, Faw et al. [18] have shown that **AdaGrad-Norm** exhibits an iteration complexity of the order of $\mathcal{O}_\ell(\epsilon^{-4})$, with high probability.

The noise adaptivity of **Adagrad** under the stronger affine condition

$$\|g_k - G_k\|^2 \stackrel{as}{\leq} \kappa_{1,\text{affine}} + \kappa_{2,\text{affine}} \|G_k\|^2 \quad (3)$$

(denoted Affine-* in Table 1) has been proved in [1]. The same complexity bound for Adagrad-Norm is obtained in [34] under the affine condition (2).

Very recently, Jiang, Maladkar and Mokhtari in [23] performed the convergence analysis in 1-norm and proved $\mathcal{O}_\ell(\epsilon^{-2})$ iteration complexity, assuming coordinate-wise bounded variance and coordinate-wise Lipschitz continuity of the gradient.

Finally Hong and Lin revisited the convergence of Adagrad in the recent paper [22], assuming a relaxed version of the condition (3) given by

$$\|g_k - G_k\|^2 \stackrel{as}{\leq} \kappa_{1,\text{affine}} + \kappa_{2,\text{affine}} \|G_k\|^2 + \kappa_{3,\text{affine}} (f(x_k) - f_{\text{low}}), \quad (4)$$

where f_{low} is such that $f(x) \geq f_{\text{low}} \forall x \in \mathbb{R}^n$. They prove that the iteration complexity interpolates between $\mathcal{O}_\ell(\epsilon^{-4})$ to $\mathcal{O}_\ell(\epsilon^{-2})$ depending on the magnitude of $\kappa_{1,\text{affine}}$ and $\kappa_{3,\text{affine}}$.

All the previous mentioned results have been obtained assuming the Lipschitz continuity of the gradient. The papers [17, 34, 22] provide results for the class of (L_0, L_1) -smooth functions² and assume that the constant L_1 is known. In [17] the results in [18] have been extended to this latter class of functions and it has been proved that the method has a complexity of the order of $\mathcal{O}_\ell(\epsilon^{-4})$, assuming (2) with $\kappa_{2,\text{affine}} < 1$. Similar results are given in [22] for Adagrad with momentum under the Assumption (4). Stronger results are obtained in [34] for AdaGrad-Norm assuming (2) where it is proved that the iteration complexity interpolates between $\mathcal{O}_\ell(\epsilon^{-4})$ to $\mathcal{O}_\ell(\epsilon^{-2})$ depending on the magnitude of $\kappa_{1,\text{affine}}$.

To conclude this brief survey, we also note that Liu et al. [29] have obtained an $\mathcal{O}(\epsilon^{-2})$ convergence result for the objective-function gap produced by the stochastic Adagrad-Norm when applied to γ -quasar convex L -smooth functions with a sub-Weibull assumption on the gradient error. This condition subsumes the sub-Gaussian case, but the result does not apply to general L -smooth nonconvex functions.

1.2 Summary of contributions

In view of all contributions discussed above, we summarize our contributions as follows.

1. We propose an extension of Adagrad which is capable of using (possibly very approximate) second-order information whenever available.
2. This algorithm is also suitable for problems involving bounds on the variables.
3. We analyze the convergence and probabilistic complexity of this algorithm, showing that it solves the approximate minimization problem with optimal complexity under a new directional condition on the gradient error, but without assuming bounded or unbiased gradient oracles nor knowledge of the problem's Lipschitz constant.
4. Taking a more general perspective, we discuss the relation between the classical optimality measure for bound-constrained problems induced by the projected gradient in the stochastic case, and have shown that, in general, the unbiased nature of the gradient oracle is not sufficient to ensure convergence on the true problem with the optimal rate, even if such a convergence occurs for the approximate one.
5. We furthermore show that our optimal complexity result for the approximated problem extends to the true problem (i.e. using exact gradients) if another condition on the gradient noise holds or if the gradient oracle is unbiased and the problem is unconstrained. We also describe the behaviour of the algorithm if none of these conditions hold.
6. As far as the authors are aware, this is the first analysis of a stochastic projected gradient method. It is also the first convergence-rate result for the stochastic Adagrad applied to general nonconvex problems that does not involve a logarithmic factor in k . This is admittedly

²A function is (L_0, L_1) -smooth if there exist a constant $L_0 \geq 0$ and $L_1 \geq 0$ and $\delta > 0$ such that for all $x, y \in \mathbb{R}^n$ with $\|x - y\| \leq \delta$, one has that $\|\nabla_x^1 f(x) - \nabla_x^1 f(y)\| \leq (L_0 + L_1 \|\nabla_x^1 f(y)\|) \|x - y\|$.

a marginal improvement of the result itself, but suggests a sharper proof technique. As it turns out, the proof is also significantly simpler than those available in the literature for the case where gradients are not uniformly bounded.

Section 2 presents the new algorithm and discusses its features, while convergence analysis for the approximate problem is described in Section 3. Section 4 discusses the relation between approximate and true optimality measures in the bound-constrained case and its application to the new algorithm. A brief conclusion and some perspectives are finally outlined in Section 5.

2 The algorithm

Our proposed algorithm, called ADAGB2, is presented on this page. Its inputs consist in a starting point x_{ini} and values l and u for the lower and upper bounds on the variables.

Algorithm 2.1: ADAGB2(x_{ini}, l, u)

Step 0: Initialization: The constants $\varsigma, \tau \in (0, 1]$ and $\kappa_s \geq 1$ are given.

Set $x_0 = P_{\mathcal{F}}(x_{\text{ini}})$, $k = 0$ and $w_{-1,i} = \varsigma$ for $i \in \{1, \dots, n\}$.

Step 1: First-order step: Compute $g_k = g(x_k, \xi)$ a random approximation of G_k , set

$$d_k \stackrel{\text{def}}{=} P_{\mathcal{F}}(x_k - g_k) - x_k, \quad (5)$$

$$w_{k,i} = \sqrt{w_{k-1,i}^2 + d_{k,i}^2} \quad \text{and} \quad \Delta_{k,i} = \frac{|d_{k,i}|}{w_{k,i}} \quad \text{for } i \in \{1, \dots, n\}. \quad (6)$$

$$\mathcal{B}_k = \{x \in \mathbb{R}^n \mid |x_i - x_{k,i}| \leq \Delta_{k,i} \text{ for } i \in \{1, \dots, n\}\} \quad (7)$$

and

$$s_k^L = P_{\mathcal{F} \cap \mathcal{B}_k}(x_k - g_k) - x_k. \quad (8)$$

Step 2: Second-order step: Choose B_k a symmetric approximation of H_k and compute

$$s_k^Q = \gamma_k s_k^L \quad \text{where} \quad \gamma_k = \begin{cases} \min \left[1, \frac{-g_k^T s_k^L}{(s_k^L)^T B_k s_k^L} \right] & \text{if } (s_k^L)^T B_k s_k^L > 0 \\ 1 & \text{otherwise.} \end{cases} \quad (9)$$

Then select s_k such that, for all $i \in \{1, \dots, n\}$,

$$x_k + s_k \in \mathcal{F}, \quad |s_{k,i}| \leq \kappa_s \Delta_{k,i} \quad \text{and} \quad g_k^T s_k + \frac{1}{2} s_k^T B_k s_k \leq \tau \left(g_k^T s_k^Q + \frac{1}{2} (s_k^Q)^T B_k s_k^Q \right). \quad (10)$$

Step 4: Loop: Set $x_{k+1} = x_k + s_k$, increment k by one and go to Step 1.

The statement of the algorithm suggests the following comments.

1. We first note that the mechanism of the algorithm ensures that all iterates remain feasible, that is $x_k \in \mathcal{F}$ for all $k \geq 0$.
2. The projections $P_{\mathcal{F}}$ and $P_{\mathcal{F} \cap \mathcal{B}_k}$ occurring in (5) and (8) are extremely cheap to compute component-wise, as, for any vector $y \in \mathbb{R}^n$ and $i \in \{1, \dots, n\}$,

$$[P_{\mathcal{F}}(y)]_i = \max[l_i, \min[y_i, u_i]] \quad (11)$$

and

$$[P_{\mathcal{F} \cap \mathcal{B}_k}(y)]_i = \max[l_i, x_{k,i} - \Delta_{k,i}, \min[y_i, x_{k,i} + \Delta_{k,i}, u_i]]. \quad (12)$$

3. When the choice $B_k = 0$ is made for all k , (9) gives that $s_k^Q = s_k^L$ and the choices $s_k = s_k^L$ or

$$s_{k,i} = -\text{sign}(g_{k,i})\Delta_{k,i}, \quad (i \in \{1, \dots, n\}) \quad (13)$$

are always acceptable for (10), in which case ADAGB2 is a purely first-order algorithm. If, in addition, the problem is unconstrained in that $l_i = -\infty$ and $u_i = +\infty$ for all i , then $d_k = -g_k$ and thus ADAGB2 with (13) is nothing but the standard Adagrad algorithm. In short,

$$\text{ADAGB2} = \text{Adagrad} + 2\text{nd order} + \text{bound constraints},$$

the last two items being of course optional. As a consequence, the stochastic complexity theory described below applies to Adagrad without any modification.

4. The uniform bound specified by AS.3 below is the only restriction made on B_k (beyond symmetry). This allows for a wide range of deterministic or random approximations, such as (sampled) Barzilai-Borwein, safeguarded (limited-memory) quasi-Newton, random sketching or finite-difference approximations of $(s_k^L)^T B_k s_k^L$ using one evaluation of the gradient oracle.
5. The shifted quadratic model of the objective function given by $g_k^T s_k + \frac{1}{2} s_k^T B s_k$ is minimized by $x_k + s_k^Q$ along the intersection of the first-order direction s_k^L with the trust region \mathcal{B}_k . In the vocabulary of trust-region methods, it can therefore be interpreted as the ‘‘Cauchy point’’ at iteration k (see [10, Sections 6.3 and 12.2.1]). Since $s_k = s_k^Q$ satisfies (10), an improved second-order step s_k beyond $x_k + s_k^Q$ is possible (for instance using a truncated conjugate-gradient or Lanczos method [33, 31] or [10, Section 7.5]), but this remains optional (as is the computation of s_k^Q itself) when access to second-order information is expensive.

3 Convergence analysis

The convergence theory we are about to describe is based on the following assumptions.

AS.1: *The function f is twice continuously differentiable and the feasible region \mathcal{F} is not empty.*

AS.2: *There exists a constant $L \geq 0$ such that for all $x, y \in \mathbb{R}^n$*

$$\|\nabla_x^1 f(x) - \nabla_x^1 f(y)\| \leq L\|x - y\|.$$

We stress that, although the existence of L is assumed, the knowledge of its value is *not* needed to run the algorithm.

AS.3: *There exists a constant $\kappa_B \geq 1$ such that $\|B_k\| \leq \kappa_B$ for all $k \geq 0$.*

AS.4: *The objective function is bounded below on the feasible domain, that is there exists a constant $f_{\text{low}} < f(x_0)$ such that $f(x) \geq f_{\text{low}}$ for every $x \in \mathcal{F}$.*

In order to state our remaining assumptions, we define the event

$$\mathcal{A} = \{\|d_0\|^2 \geq \varsigma\}. \quad (14)$$

This event occurs or does not occur at iteration 0, i.e. at the very beginning of a realization of the ADAGB2 algorithm. The convergence theory which follows is dependent of the (in practice extremely likely) occurrence of \mathcal{A} and our subsequent stochastic assumptions are therefore conditioned by this event. They will be formally specified by considering, at iteration j , expectations conditioned by the past iterations and by \mathcal{A} , which will be denoted by the symbol $\mathbb{E}_j^{\mathcal{A}}[\cdot]$. Note that, because \mathcal{A} is measurable for all iterations after the zero-th one, $\mathbb{E}_k^{\mathcal{A}}[\cdot] = \mathbb{E}_k[\cdot]$ whenever $k > 0$.

AS.5: *There exists a constant $\kappa_{Gg} > 0$ such that*

$$\mathbb{E}_k^{\mathcal{A}}[|(G_k - g_k)^T s_k|] \leq \kappa_{Gg}^2 \mathbb{E}_k^{\mathcal{A}}[\|s_k\|^2] \quad (15)$$

for all $k \geq 0$.

This condition (which we have called "new1" in Table 1) can be interpreted as a "directional root mean square" condition along the direction s_k (see the comments after Theorem 4.2 below). As far as the authors are aware, this condition is necessary if a Cauchy point is introduced to take second-order information into account because γ_k then becomes a random variable.

The following "linear descent" lemma is a variant of [20, Lemma 2.1] (which, in the unconstrained deterministic context, uses a different optimality measure and a different definition of s_k^L).

Lemma 3.1 Suppose that AS.3 and AS.5 hold. Then, for $j \geq 0$,

$$\mathbb{E}_j^{\mathcal{A}}[G_j^T s_j] \leq -\frac{\tau\varsigma^2}{2\kappa_B} \mathbb{E}_j^{\mathcal{A}}[d_j^T \Delta_j] + \kappa_s^2 \left(\frac{1}{2}\kappa_B + \kappa_{Gg}^2 \right) \mathbb{E}_j^{\mathcal{A}}[\|\Delta_j\|^2]. \quad (16)$$

Proof. Consider any component $i \in \{1, \dots, n\}$. We first note that (5) and the contractive nature of the projection $P_{\mathcal{F}}$ ensure that $|g_{j,i}| \geq |d_{j,i}|$. Moreover, (5), (6), (7) and (8) implies that either $|s_{j,i}^L| = |d_{j,i}|/w_{j,i}$ or $|s_{j,i}^L| = |d_{j,i}|$. In the first case, we may deduce that

$$|g_{j,i} s_{j,i}^L| \geq \frac{d_{j,i}^2}{w_{j,i}} = w_{j,i} \frac{d_{j,i}^2}{w_{j,i}^2} \geq \varsigma (s_{j,i}^L)^2,$$

while, in the latter case,

$$|g_{j,i} s_{j,i}^L| \geq d_{j,i}^2 = (s_{j,i}^L)^2 = w_{j,i} \frac{d_{j,i}^2}{w_{j,i}} \geq \varsigma \frac{d_{j,i}^2}{w_{j,i}}.$$

Combining the two cases and remembering that $\varsigma \in (0, 1]$, we see that, for $i \in \{1, \dots, n\}$,

$$|g_{j,i} s_{j,i}^L| \geq \varsigma \frac{d_{j,i}^2}{w_{j,i}} \quad \text{and} \quad |g_{j,i} s_{j,i}^L| \geq \varsigma (s_{j,i}^L)^2.$$

Summing over all components $i \in \{1, \dots, n\}$ and using the fact that, by construction, $g_{j,i} s_{j,i}^L < 0$ for all i then gives that

$$g_j^T s_j^L \leq -\varsigma \sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}} \quad \text{and} \quad |g_j^T s_j^L| \geq \varsigma \|s_j^L\|^2. \quad (17)$$

We now consider the quadratic model and suppose first that $\gamma_j < 1$. Thus $(s_j^L)^T B_j s_j^L > 0$. We then deduce from (9), AS.3 and (17) that

$$g_j^T s_j^Q + \frac{1}{2}(s_j^Q)^T B_j s_j^Q = -\frac{(g_j^T s_j^L)^2}{2(s_j^L)^T B_j s_j^L} \leq -\frac{\varsigma}{2\kappa_B} |g_j^T s_j^L| \leq -\frac{\varsigma^2}{2\kappa_B} \sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}}. \quad (18)$$

If now $\gamma_j = 1$, then (9) and (17) give that

$$g_j^T s_j^Q + \frac{1}{2}(s_j^Q)^T B_j s_j^Q = g_j^T s_j^L + \frac{1}{2}(s_j^L)^T B_j s_j^L \leq \frac{1}{2} g_j^T s_j^L < -\frac{\varsigma}{2} \sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}}$$

and (18) then again follows from the bounds $\kappa_B \geq 1$ and $\varsigma \leq 1$. Thus, successively using the third part of (10), AS.3, the second part of (10) and (18),

$$\begin{aligned} g_j^T s_j &= g_j^T s_j + \frac{1}{2} s_j^T B_j s_j - \frac{1}{2} s_j^T B_j s_j \\ &\leq \tau \left(g_j^T s_j^Q + \frac{1}{2} (s_j^Q)^T B_j s_j^Q \right) + \frac{1}{2} \kappa_B \|s_j\|^2 \\ &\leq \tau \left(g_j^T s_j^Q + \frac{1}{2} \frac{1}{2} (s_j^Q)^T B_j s_j^Q \right) + \frac{1}{2} \kappa_s^2 \kappa_B \|\Delta_j\|^2 \\ &\leq -\frac{\tau \varsigma^2}{2\kappa_B} \sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}} + \frac{1}{2} \kappa_s^2 \kappa_B \|\Delta_j\|^2, \end{aligned}$$

which, with the second part of (6), gives that

$$g_j^T s_j \leq -\frac{\tau \varsigma^2}{2\kappa_B} d_j^T \Delta_j + \frac{1}{2} \kappa_s^2 \kappa_B \|\Delta_j\|^2. \quad (19)$$

But, using AS.5 and the second part of (10),

$$\begin{aligned} \mathbb{E}_j^A [G_j^T s_j] &= \mathbb{E}_j^A [g_j^T s_j] + \mathbb{E}_j^A [(G_j - g_j)^T s_j] \\ &\leq \mathbb{E}_j^A [g_j^T s_j] + |\mathbb{E}_j^A [(G_j - g_j)^T s_j]| \\ &\leq \mathbb{E}_j^A [g_j^T s_j] + \kappa_{Gg}^2 |\mathbb{E}_j^A [\|s_j\|^2]| \\ &\leq \mathbb{E}_j^A [g_j^T s_j] + \kappa_{Gg}^2 \kappa_s^2 \mathbb{E}_j^A [\|\Delta_j\|^2] \end{aligned}$$

and (16) follows from (19). \square

Lemma 3.1 is crucial for the proof of our complexity bounds, which also depends on two known technical results. We restate them for completeness.

Lemma 3.2 Let $\{a_j\}_{j \geq 0}$ be a sequence of non-negative numbers and let $b_j = \sum_{i=0}^j a_i$. Then

$$\sum_{j=0}^k \frac{a_j}{\varsigma + b_j} \leq \log \left(1 + \frac{1}{\varsigma} b_k \right).$$

Proof. See [37]. \square

Lemma 3.3 Suppose that

$$\gamma_1 u \leq \gamma_2 \log(u) \quad (20)$$

for some $\gamma_1, \gamma_2 > 0$ with $\gamma_2 > 3\gamma_1$. Then

$$0 < u \leq -\frac{\gamma_2}{\gamma_1} W_{-1} \left(-\frac{\gamma_1}{\gamma_2} \right) \leq \frac{\gamma_2}{\gamma_1} \left[\log \left(\frac{\gamma_2}{\gamma_1} \right) + \sqrt{2 \left(\log \left(\frac{\gamma_2}{\gamma_1} \right) - 1 \right)} \right] \quad (21)$$

where W_{-1} is the second branch of the Lambert function [11].

Proof. [Extracted from the proof of [21, Theorem 3.2]].

Let $\psi(u) \stackrel{\text{def}}{=} \gamma_1 u - \gamma_2 \log(u)$. Since $\gamma_2 > 3\gamma_1$, the equation $\psi(u) = 0$ admits two roots $u_1 \leq u_2$ and (20) holds for $u \in [u_1, u_2]$. The definition of u_2 then gives that

$$\log(u_2) - \frac{\gamma_1}{\gamma_2} u_2 = 0$$

that is

$$u_2 e^{-\frac{\gamma_1}{\gamma_2} u_2} = 1.$$

If $z = -\frac{\gamma_1}{\gamma_2} u_2$, we thus obtain that

$$z e^z = -\frac{\gamma_1}{\gamma_2}$$

Hence $z = W_{-1}(-\frac{\gamma_1}{\gamma_2}) < 0$, where W_{-1} is the second branch of the Lambert function on $[-\frac{1}{e}, 0)$.

As $-\frac{\gamma_1}{\gamma_2} \geq -\frac{1}{3}$, z is well defined and therefore

$$u_2 = -\frac{\gamma_2}{\gamma_1} z = -\frac{\gamma_2}{\gamma_1} W_{-1} \left(-\frac{\gamma_1}{\gamma_2} \right) > 0,$$

proving the first inequality of (21). The second inequality is obtained by exploiting an upper bound on the value of the Lambert function [8, Theorem 1] which states that, for $x > 0$,

$$|W_{-1}(-e^{-x-1})| \leq 1 + \sqrt{2x} + x.$$

Taking $x = \log \left(\frac{\gamma_2}{\gamma_1} \right) - 1 > 0$ then gives the desired result. \square

Our first complexity result now consider the global rate of convergence of $\|d_k\|$ to zero. It is expressed using the expectation conditioned by \mathcal{A} , which we denote by the symbol $\mathbb{E}^{\mathcal{A}}[\cdot]$.

Theorem 3.4 Suppose that AS.1–AS.5 hold and that the ADAGB2 algorithm is applied to problem (1). Then

$$\mathbb{E}^{\mathcal{A}} \left[\text{average}_{j \in \{0, \dots, k\}} \|d_j\| \right] \leq \frac{\kappa_{\text{conv}}}{\sqrt{k+1}}, \quad (22)$$

with

$$\kappa_{\text{conv}} = \sqrt{\frac{\zeta}{2}} \kappa_W \left| W_{-1} \left(-\frac{1}{\kappa_W} \right) \right| \leq \sqrt{\frac{\zeta}{2}} \kappa_W \left| \log(\kappa_W) + \sqrt{2(\log(\kappa_W) - 1)} \right|, \quad (23)$$

where W_{-1} is the second branch of the Lambert function, $\Gamma_0 \stackrel{\text{def}}{=} f(x_0) - f_{\text{low}}$ and

$$\kappa_W \stackrel{\text{def}}{=} \frac{8\kappa_B}{\tau\zeta^2\sqrt{\zeta}} \max[1, \Gamma_0] \max \left[3, \frac{n\kappa_s^2}{\Gamma_0} (2\kappa_{Gg}^2 + \kappa_B + L) \right]. \quad (24)$$

Proof. Consider an iteration j and first note that

$$\|s_j\|^2 \leq \kappa_s^2 \|\Delta_j\|^2 = \kappa_s^2 \sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}^2}$$

because of the second part of (10) and (6). Lemma 3.1, AS.1 and AS.2 then give that

$$\begin{aligned} \mathbb{E}_j^{\mathcal{A}}[f(x_{j+1})] &\leq f(x_j) + \mathbb{E}_j^{\mathcal{A}}[G_j^T s_j] + \frac{L}{2} \mathbb{E}_j^{\mathcal{A}}[\|s_j\|^2] \\ &\leq f(x_j) - \frac{\tau\varsigma^2}{2\kappa_B} \mathbb{E}_j^{\mathcal{A}} \left[\sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}} \right] + \kappa_s^2 (\kappa_{Gg}^2 + \frac{1}{2}\kappa_B + \frac{1}{2}L) \mathbb{E}_j^{\mathcal{A}} \left[\sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}^2} \right]. \end{aligned} \quad (25)$$

Defining $\kappa_* = \kappa_s^2 (\kappa_{Gg}^2 + \frac{1}{2}\kappa_B + \frac{1}{2}L)$, taking the expectation conditional to \mathcal{A} and using the tower property gives that

$$\mathbb{E}^{\mathcal{A}}[f(x_{j+1})] \leq \mathbb{E}^{\mathcal{A}}[f(x_j)] - \frac{\tau\varsigma^2}{2\kappa_B} \mathbb{E}^{\mathcal{A}} \left[\sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}} \right] + \kappa_* \mathbb{E}^{\mathcal{A}} \left[\sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}^2} \right],$$

and therefore, summing for $j \in \{0, \dots, k\}$ for k fixed, that

$$\frac{\tau\varsigma^2}{2\kappa_B} \mathbb{E}^{\mathcal{A}} \left[\sum_{j=0}^k \sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}} \right] \leq f(x_0) - f_{\text{low}} + \kappa_* \mathbb{E}^{\mathcal{A}} \left[\sum_{j=0}^k \sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}^2} \right]. \quad (26)$$

Now, for every realization such that \mathcal{A} occurs and using the non-decreasing nature of $w_{j,i}$ as a function of j ,

$$\sum_{j=0}^k \sum_{i=1}^n \frac{d_{j,i}^2}{\max_{\ell \in \{1, \dots, n\}} w_{k,\ell}} \leq \sum_{j=0}^k \sum_{i=1}^n \frac{d_{j,i}^2}{w_{j,i}}.$$

We may now apply Lemma 3.2 to derive that, for each $i \in \{1, \dots, n\}$,

$$\begin{aligned} \mathbb{E}^{\mathcal{A}} \left[\sum_{i=1}^n \sum_{j=0}^k \frac{d_{j,i}^2}{w_{j,i}^2} \right] &= \mathbb{E}^{\mathcal{A}} \left[\sum_{i=1}^n \sum_{j=0}^k \frac{d_{j,i}^2}{\varsigma + \sum_{i=0}^j d_{j,i}^2} \right] \\ &\leq n \mathbb{E}^{\mathcal{A}} \left[\max_{i \in \{1, \dots, n\}} \log \left(1 + \frac{1}{\varsigma} \sum_{j=0}^k d_{j,i}^2 \right) \right] \\ &\leq n \mathbb{E}^{\mathcal{A}} \left[\log \left(1 + \frac{1}{\varsigma} \sum_{i=1}^n \sum_{j=0}^k d_{j,i}^2 \right) \right] \end{aligned}$$

and hence, using AS.4 to define $\Gamma_0 = f(x_0) - f_{\text{low}} > 0$, (27) gives

$$\frac{\tau\varsigma^2}{2\kappa_B} \mathbb{E}^{\mathcal{A}} \left[\sum_{j=0}^k \sum_{i=1}^n \frac{d_{j,i}^2}{\max_{\ell \in \{1, \dots, n\}} w_{k,\ell}} \right] \leq \Gamma_0 + n\kappa_* \mathbb{E}^{\mathcal{A}} \left[\log \left(1 + \frac{1}{\varsigma} \sum_{j=0}^k \|d_j\|^2 \right) \right] \quad (27)$$

Now for any realization for which \mathcal{A} occurs,

$$\sum_{j=0}^k \|d_j\|^2 \geq \|d_0\|^2 \geq \varsigma \quad (28)$$

and thus

$$\max_{\ell \in \{1, \dots, n\}} w_{k,\ell} < \sqrt{\varsigma + \sum_{j=0}^k \|d_j\|^2} \leq \sqrt{2 \sum_{j=0}^k \|d_j\|^2} \quad (29)$$

and

$$2 \leq 1 + \frac{1}{\varsigma} \sum_{j=0}^k \|d_j\|^2 \leq \frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2,$$

so that, using the monotonicity of the logarithm,

$$0 < \log(\sqrt{2}) \leq \log \left(\mathbb{E}^{\mathcal{A}} \left[\sqrt{\frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2} \right] \right). \quad (30)$$

Inserting (29) in (27) then yields that

$$\begin{aligned} \frac{\tau\varsigma^2}{2\sqrt{2}\kappa_B} \mathbb{E}^{\mathcal{A}} \left[\frac{\sum_{j=0}^k \|d_j\|^2}{\sqrt{\sum_{j=0}^k \|d_j\|^2}} \right] &\leq \Gamma_0 + n\kappa_* \mathbb{E}^{\mathcal{A}} \left[\log \left(\frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2 \right) \right] \\ &\leq \Gamma_0 + 2n\kappa_* \mathbb{E}^{\mathcal{A}} \left[\log \left(\sqrt{\frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2} \right) \right] \end{aligned}$$

and thus, using the concavity of the logarithm, that

$$\frac{\tau\varsigma^2\sqrt{\varsigma}}{4\kappa_B} \left(\mathbb{E}^{\mathcal{A}} \left[\sqrt{\frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2} \right] \right) \leq \Gamma_0 + 2n\kappa_* \log \left(\mathbb{E}^{\mathcal{A}} \left[\sqrt{\frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2} \right] \right). \quad (31)$$

Dividing both sides by Γ_0 and using (30) then gives that

$$\begin{aligned} \frac{\tau\varsigma^2\sqrt{\varsigma}}{4\kappa_B \max[1, \Gamma_0]} \left(\mathbb{E}^{\mathcal{A}} \left[\sqrt{\frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2} \right] \right) &\leq 1 + \frac{2n\kappa_*}{\Gamma_0} \log \left(\mathbb{E}^{\mathcal{A}} \left[\sqrt{\frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2} \right] \right) \\ &\leq \left(\frac{1}{\log(\sqrt{2})} + \frac{2n\kappa_*}{\Gamma_0} \right) \log \left(\mathbb{E}^{\mathcal{A}} \left[\sqrt{\frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2} \right] \right) \\ &\leq 2 \max \left[3, \frac{2n\kappa_*}{\Gamma_0} \right] \log \left(\mathbb{E}^{\mathcal{A}} \left[\sqrt{\frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2} \right] \right). \end{aligned} \quad (32)$$

We may now apply Lemma 3.3 since (32) is identical to (20) with

$$\gamma_1 \stackrel{\text{def}}{=} \frac{\tau\varsigma^2\sqrt{\varsigma}}{4\kappa_B \max[1, \Gamma_0]} \quad \gamma_2 \stackrel{\text{def}}{=} 2 \max \left[3, \frac{2n\kappa_*}{\Gamma_0} \right] \quad \text{and} \quad u \stackrel{\text{def}}{=} \mathbb{E}^{\mathcal{A}} \left[\sqrt{\frac{2}{\varsigma} \sum_{j=0}^k \|d_j\|^2} \right] \quad (33)$$

and $\gamma_1 \leq 1$ (because $\tau \leq 1$, $\varsigma \leq 1$ and $\kappa_B \geq 1$) ensures that $\gamma_2 > 3\gamma_1$. If we now define

$$\kappa_W = \frac{\gamma_2}{\gamma_1} = \frac{8\kappa_B}{\tau\varsigma^2\sqrt{\varsigma}} \max[1, \Gamma_0] \max \left[3, \frac{2n\kappa_*}{\Gamma_0} \right], \quad (34)$$

we obtain from (33) and the first part of (21) that

$$\mathbb{E}^{\mathcal{A}} \left[\sqrt{\sum_{j=0}^k \|d_j\|^2} \right] \leq \sqrt{\frac{\varsigma}{2}} \kappa_W \left| W_{-1} \left(-\frac{1}{\kappa_W} \right) \right|.$$

Dividing by $\sqrt{k+1}$ and using the inequality

$$\frac{1}{k+1} \sum_{j=0}^k \|d_j\| \leq \frac{1}{\sqrt{k+1}} \sqrt{\sum_{j=0}^k \|d_j\|^2}$$

finally gives

$$\mathbb{E}^{\mathcal{A}} \left[\text{average}_{j \in \{0, \dots, k\}} \|d_j\| \right] \leq \mathbb{E}^{\mathcal{A}} \left[\sqrt{\text{average}_{j \in \{0, \dots, k\}} \|d_j\|^2} \right] \leq \left(\sqrt{\frac{\varsigma}{2}} \kappa_W \left| W_{-1} \left(-\frac{1}{\kappa_W} \right) \right| \right) \cdot \frac{1}{\sqrt{k+1}}.$$

Substituting the value of κ_* in (34) and using the second inequality in (21) then concludes the proof. \square

We now briefly return to our assumption that \mathcal{A} holds to derive the above results. In fact, this assumption is only made to ensure (28) at iteration k . But (28) also shows that assuming the occurrence of \mathcal{A} may be replaced by the even weaker requirement that $\sum_{j=0}^k \|d_j\|^2 \geq \varsigma$ suffices.

In the deterministic case, Theorem 3.4 provides a bound on the complexity of solving the bound-constrained problem (1) for which we consider the standard optimality measure [10, Section 12.1] $\|\Xi_k\|$, where

$$\Xi_k \stackrel{\text{def}}{=} P_{\mathcal{F}}(x_k - G_k) - x_k$$

(if the problem is unconstrained, $\|\Xi_k\| = \|G_k\|$).

Corollary 3.5 Suppose that AS.1–AS.4 hold, that the ADAGB2 algorithm is applied to problem (1) starting from a non-critical x_0 with $\varsigma < \|d_0\|^2$ and that $g_j = G_j$ for all j . Then

$$\min_{j \in \{0, \dots, k\}} \|\Xi_j\| \leq \text{average}_{j \in \{0, \dots, k\}} \|\Xi_j\| \leq \frac{\kappa_{\text{conv}}}{\sqrt{k+1}}, \quad (35)$$

where the constant κ_{conv} is computed as in Theorem 3.4 using the value $\kappa_{Gg} = 0$.

Proof. The choice $\varsigma < \|d_0\|^2$ (which is always possible in the deterministic context) ensures that \mathcal{A} occurs and AS.5 holds because $g_j = G_j$. The result then follows from Theorem 3.4 with $\kappa_{Gg} = 0$. \square

4 First-order optimality and stochastic projected gradients with bounds

4.1 A more general framework

Things are more complicated in the stochastic case, as we examine in this section. However, as our arguments apply not only to the ADAGB2 algorithm, but also to a wider class of stochastic projected-gradient methods, we now consider solving problem (1) using the algorithmic framework given by StochProjGrad.

Moreover we will assume, in this subsection, that the expected approximate criticality measure

$$\mathbb{E}_k[\|d_k\|] = \mathbb{E}_k[\|P_{\mathcal{F}}(x_k - g_k) - x_k\|]$$

converges to zero (at least on average as in Theorem 3.4). We are now interested in what can be deduced on $\mathbb{E}[\|\Xi_k\|] = \mathbb{E}[\|P_{\mathcal{F}}(x_k - G_k) - x_k\|]$, the relevant criticality measure for problem (1) in

Algorithm 4.1: StochProjGrad (x_{ini}, l, u)

Step 0: Initialization: Set $x_0 = P_{\mathcal{F}}(x_{\text{ini}})$, $k = 0$.

Step 1: Step: Compute $g_k = g(x_k, \xi)$ a random approximation of $G_k = \nabla_x^1 f(x_k)$, and a step s_k such that $x_k + s_k \in \mathcal{F}$,

Step 4: Loop: Set $x_{k+1} = x_k + s_k$, increment k by one and go to Step 1.

the stochastic case. Ideally, one would hope that the approximate gradient's distribution ensures coherence of the measures in the sense that

$$\mathbb{E}[\|\Xi_k\|] \leq \kappa_{\text{opt}} \mathbb{E}[\|d_k\|] \quad (36)$$

for a fixed $\kappa_{\text{opt}} > 0$, in which case $\mathbb{E}[\|\Xi_k\|]$ converges to zero at the same rate as $\mathbb{E}[\|d_k\|]$. If we consider the unconstrained case ($\mathcal{F} = \mathbb{R}^n$) and assume that the gradient oracle is unbiased (i.e. $\mathbb{E}_k[g_k] = G_k$), then, using Jensen's inequality and the convexity of the norm,

$$\|\Xi_k\| = \|G_k\| = \|\mathbb{E}_k[g_k]\| \leq \mathbb{E}_k[\|g_k\|] = \mathbb{E}_k[\|d_k\|].$$

Taking the full expectation on both sides and using the law of total expectation then shows that (36) always holds with $\kappa_{\text{opt}} = 1$. The situation is more complicated if bounds are present, even if the gradient oracle is unbiased. To see this, consider the following one dimensional example, where $\mathcal{F} = [0, +\infty)$ and

$$x_k = \frac{1}{k+1} \quad \text{and} \quad g_k = \begin{cases} 1 & \text{with probability } p_k = \frac{1}{k+1} + \frac{1}{(k+1)^2} \\ 0 & \text{with probability } 1 - p_k, \end{cases}$$

for $k > 0$. Also define $G_k = \mathbb{E}_k[g_k] = p_k$ (so that the gradient oracle is unbiased). One then easily verifies that $|G_{k+1} - G_k| \leq 3|x_{k+1} - x_k|$, so that AS.2 holds with $L = 3$. These definitions give that

$$|\Xi_k| = \left| P_{[0,+\infty)} \left(\frac{1}{k+1} - p_k \right) - \frac{1}{k+1} \right| = \left| P_{[0,+\infty)} \left(\frac{-1}{(k+1)^2} \right) - \frac{1}{k+1} \right| = \frac{1}{k+1}$$

and

$$\begin{aligned} \mathbb{E}_k[\|d_k\|] &= \left| P_{[0,+\infty)} \left(\frac{1}{k+1} - 1 \right) - \frac{1}{k+1} \right| p_k + \left| P_{[0,+\infty)} \left(\frac{1}{k+1} - 0 \right) - \frac{1}{k+1} \right| (1 - p_k) \\ &= \frac{1}{k+1} p_k + 0(1 - p_k) \\ &= \frac{1}{(k+1)^2} + \frac{1}{(k+1)^3}. \end{aligned}$$

Applying now the law of total expectation, we see that

$$\lim_{k \rightarrow \infty} \frac{\mathbb{E}[\|d_k\|]}{\mathbb{E}[\|\Xi_k\|]} = 0,$$

preventing (36) (a "coherently distributed" gradient oracle) to hold for a fixed $\kappa_{\text{opt}} > 0$. We therefore conclude that, in general, *the rate of decrease of $\mathbb{E}[\|d_k\|]$ does not translate to a similar rate of decrease for $\mathbb{E}[\|\Xi_k\|]$, even for unbiased gradient oracles.*

Fortunately, the situation can be improved by strengthening the condition on the gradient accuracy, even if the gradient oracle is biased. This is the object of the next lemma.

Lemma 4.1 For each $k \geq 0$ and each $i \in \{1, \dots, n\}$, we have that

$$\mathbb{E}[\|\Xi_k\|] \leq \mathbb{E}[\|d_k\|] + \mathbb{E}[\|g_k - G_k\|]. \quad (37)$$

Moreover, if $\mathbb{E}_k[\|g_k - G_k\|] \leq \kappa_{\text{err}} \mathbb{E}_k[\|d_k\|]$ for some $\kappa_{\text{err}} \geq 0$, then

$$\mathbb{E}[\|\Xi_k\|] \leq (1 + \kappa_{\text{err}}) \mathbb{E}[\|d_k\|]. \quad (38)$$

Proof. Consider an arbitrary $i \in \{1, \dots, n\}$ and note that

$$\Xi_{k,i} = P_i(x_{k,i} - G_{k,i}) - x_{k,i}. \quad (39)$$

where $P_i = P_{[l_i, u_{l_i}]}$. Since P_i is a contractive map, we have that

$$\begin{aligned} & (P_i(x_{k,i} - G_{k,i}) - x_{k,i})^2 - (P_i(x_{k,i} - g_{k,i}) - x_{k,i})^2 \\ &= [P_i(x_{k,i} - G_{k,i}) - x_{k,i} + P_i(x_{k,i} - g_{k,i}) - x_{k,i}] \times \\ & \quad [P_i(x_{k,i} - G_{k,i}) - P_i(x_{k,i} - g_{k,i})] \\ & \leq |P_i(x_{k,i} - G_{k,i}) - x_{k,i} + P_i(x_{k,i} - g_{k,i}) - x_{k,i}| |g_{k,i} - G_{k,i}|. \end{aligned}$$

Now, again using the contractivity of P_i ,

$$\begin{aligned} & |P_i(x_{k,i} - G_{k,i}) - x_{k,i} + P_i(x_{k,i} - g_{k,i}) - x_{k,i}| \\ &= |2(P_i(x_{k,i} - g_{k,i}) - x_{k,i}) \\ & \quad + (P_i(x_{k,i} - G_{k,i}) - x_{k,i}) - (P_i(x_{k,i} - g_{k,i}) - x_{k,i})| \\ & \leq 2|P_i(x_{k,i} - g_{k,i}) - x_{k,i}| + |P_i(x_{k,i} - G_{k,i}) - P_i(x_{k,i} - g_{k,i})| \\ & \leq 2|P_i(x_{k,i} - g_{k,i}) - x_{k,i}| + |g_{k,i} - G_{k,i}| \end{aligned}$$

and thus

$$\begin{aligned} & (P_i(x_{k,i} - G_{k,i}) - x_{k,i})^2 - (P_i(x_{k,i} - g_{k,i}) - x_{k,i})^2 \\ & \leq 2|P_i(x_{k,i} - g_{k,i}) - x_{k,i}| |g_{k,i} - G_{k,i}| + |g_{k,i} - G_{k,i}|^2, \end{aligned}$$

which, using (5) and (39), gives that

$$\Xi_{k,i}^2 \leq d_{k,i}^2 + 2|d_{k,i}| |g_{k,i} - G_{k,i}| + |g_{k,i} - G_{k,i}|^2$$

and therefore, summing for $i \in \{1, \dots, n\}$ and using the Cauchy-Schwartz inequality, that

$$\begin{aligned} \|\Xi_k\|^2 &= \|d_k\|^2 + 2|d_k|^T |g_k - G_k| + \|g_k - G_k\|^2 \\ &\leq \|d_k\|^2 + 2\|d_k\| \|g_k - G_k\| + \|g_k - G_k\|^2 \\ &= \left(\|d_k\| + \|g_k - G_k\| \right)^2. \end{aligned}$$

This yields that

$$\|\Xi_k\| \leq \|d_k\| + \|g_k - G_k\|. \quad (40)$$

and (37) follows by taking conditional expectations on both sides. Moreover, if $\mathbb{E}_k[\|g_k - G_k\|] \leq \kappa_{\text{err}} \mathbb{E}_k[\|d_k\|]$, (38) results from (40), and the tower property. \square

Thus we see from (37) that the convergence of $\mathbb{E}[\|\Xi_k\|]$ to zero is still guaranteed if $\mathbb{E}[\|g_k - G_k\|]$ converges to zero along with $\mathbb{E}[\|d_k\|]$. Moreover, if it does so at the same speed, the rate of decrease of $\mathbb{E}[\|d_k\|]$ dominates, as shown by (38).

This result also indicates what can happen if nothing is known on the error in the gradient. The optimality measure of the approximate problem goes to zero (as in Theorem 3.4) but the true measure could remain at a level given by the second term of (37). To interpret this term, observe that, using Jensen's inequality and the concavity of the square root,

$$\mathbb{E}_k[\|g_k - G_k\|] = \mathbb{E}_k\left[\sqrt{\|g_k - G_k\|^2}\right] \leq \sqrt{\mathbb{E}_k[\|g_k - G_k\|^2]} = \text{RMSE}_k \quad (41)$$

for all $k \geq 0$, where RMSE_k is the conditional root mean square error (RMSE) of the gradient oracle at iteration k . Thus, by the law of total expectation, $\mathbb{E}[\|g_k - G_k\|] \leq \mathbb{E}[\text{RMSE}_k]$, and $\limsup_{k \rightarrow \infty} \mathbb{E}[\text{RMSE}_k]$ gives an asymptotic upper bound on the criticality default (the deviation of the criticality measure from zero).

4.2 Application to the ADAGB2 algorithm

The discussion of the previous subsection (where we may obviously condition every expectation to \mathcal{A}) finally allows us to rephrase Theorem 3.4 to cover convergence of the ADAGB2 algorithm on problem (1) in three progressively more general scenarii (condition (44) below was called "new2" in Table 1).

Theorem 4.2 Suppose that AS.1–AS.5 hold and that the ADAGB2 algorithm is applied to problem (1). Then

$$\text{average}_{j \in \{0, \dots, k\}} \mathbb{E}^{\mathcal{A}}[\|\Xi_j\|] \leq \frac{\kappa_{\text{conv}}}{\sqrt{k+1}} + \kappa_2 \text{average}_{j \in \{0, \dots, k\}} \mathbb{E}^{\mathcal{A}}[\|g_k - G_k\|] \quad (42)$$

with

$$\kappa_{\text{conv}} = \kappa_1 \sqrt{\frac{\zeta}{2}} \kappa_W \left| W_{-1} \left(-\frac{1}{\kappa_W} \right) \right| \leq \kappa_1 \sqrt{\frac{\zeta}{2}} \kappa_W \left| \log(\kappa_W) + \sqrt{2(\log(\kappa_W) - 1)} \right|, \quad (43)$$

where W_{-1} is the second branch of the Lambert function, $\Gamma_0 \stackrel{\text{def}}{=} f(x_0) - f_{\text{low}}$, κ_W is defined in (24) and

coherently distributed: $\kappa_1 = \kappa_{\text{opt}}$ and $\kappa_2 = 0$ if

$$\mathbb{E}^{\mathcal{A}}[\|\Xi_k\|] \leq \kappa_{\text{opt}} \mathbb{E}^{\mathcal{A}}[\|d_k\|]$$

for some constant $\kappa_{\text{opt}} > 0$ and all $k \geq 0$;

controlled error: $\kappa_1 = 1 + \kappa_{\text{err}}$ and $\kappa_2 = 0$ if

$$\mathbb{E}_k^{\mathcal{A}}[\|g_k - G_k\|] \leq \kappa_{\text{err}} \mathbb{E}_k^{\mathcal{A}}[\|d_k\|] \quad (44)$$

for some constant $\kappa_{\text{err}} > 0$ and all $k \geq 0$;

general: $\kappa_1 = \kappa_2 = 1$ otherwise.

This theorem gives the desired fast convergence to zero of the first-order optimality measure associated with problem (1) in the coherently distributed and controlled error cases. In particular, if the problem is unconstrained and the gradient oracle is unbiased, then the coherently distributed case applies with $\kappa_{\text{opt}} = 1$ and we recover the $\mathcal{O}(\epsilon^{-2})$ complexity bound for Adagrad obtained in the

references mentioned in the introduction *only assuming the directional condition given by AS.5*. The same is true in the bound-constrained case if controlled error is assumed.

Theorem 4.2 and the discussion of the previous section also indicate that the true measure could remain at a level given by $\limsup_{k \rightarrow \infty} \beta_k$ where $\beta_k = \text{average}_{j \in \{0, \dots, k\}} \mathbb{E}^{\mathcal{A}}[\|g_k - G_k\|]$ is bounded above by the average $\text{RMSE}_k^{\mathcal{A}}$ of the gradient oracle, the average being taken on the first k iterations. If β_k is bounded for large k by a (hopefully) small positive constant, then the optimality measure is only constrained to fall below this constant. However, convergence to true optimality must still happen if the sequence $\{\beta_k\}$ converges to zero³, albeit possibly at a rate slower than $\mathcal{O}(1/\sqrt{k+1})$.

Moreover, reformulating the "new2" condition (44) using the bound (41) allows us to state the following direct consequence of Theorem 4.2.

Corollary 4.3 Suppose that AS.1–AS.5 hold and that the ADAGB2 algorithm is applied to problem (1). Suppose also that there exists a constant $\kappa_{\text{err}} \geq 0$ such that, for all $k \geq 0$,

$$\text{RMSE}_k^{\mathcal{A}} \leq \kappa_{\text{err}} \mathbb{E}^{\mathcal{A}}[\|d_k\|]$$

where $\text{RMSE}_k^{\mathcal{A}}$ is the conditional root mean square error defined by the last equality of (41). Then

$$\text{average}_{j \in \{0, \dots, k\}} \mathbb{E}^{\mathcal{A}}[\|\Xi_j\|] \leq \frac{\kappa_{\text{conv}}}{\sqrt{k+1}},$$

where κ_{conv} defined by (43).

We finally state a complexity result in probability simply derived from Theorem 4.2 for its first two scenarii.

Corollary 4.4 Suppose that the conditions of Theorem 4.2 hold for the coherently distributed or controlled error cases. Then, given $\delta \in (1 - p_{\mathcal{A}}, 1)$ where $p_{\mathcal{A}}$ is the probability of occurrence of the event \mathcal{A} , one has that

$$\mathbb{P} \left[\min_{j \in \{0, \dots, k\}} \|\Xi_j\| \leq \epsilon \right] \geq 1 - \delta \quad \text{for } k \geq \left(\frac{p_{\mathcal{A}} \kappa_{\text{conv}}}{(p_{\mathcal{A}} - (1 - \delta)) \epsilon} \right)^2,$$

where κ_{conv} is given by (43).

Proof. Using Markov's inequality and (42) with $\kappa_{\text{bias}} = 0$, we have that

$$\mathbb{P}^{\mathcal{A}} \left[\min_{j \in \{0, \dots, k\}} \|\Xi_j\| \leq \epsilon \right] \geq \mathbb{P}^{\mathcal{A}} \left[\frac{1}{k} \sum_{j=0}^k \|\Xi_j\| \leq \epsilon \right] \geq 1 - \frac{1}{\epsilon} \mathbb{E}_k^{\mathcal{A}} \left[\frac{1}{k} \sum_{j=0}^k \|\Xi_j\| \right] \geq 1 - \frac{\kappa_{\text{conv}}}{\epsilon \sqrt{k+1}},$$

and thus, given the definition of $p_{\mathcal{A}}$, we deduce that

$$\mathbb{P} \left[\min_{j \in \{1, \dots, k\}} \|\Xi_j\| \leq \epsilon \right] \geq p_{\mathcal{A}} \mathbb{P}^{\mathcal{A}} \left[\min_{j \in \{1, \dots, k\}} \|\Xi_j\| \leq \epsilon \right] \geq p_{\mathcal{A}} \left(1 - \frac{\kappa_{\text{conv}}}{\epsilon \sqrt{k+1}} \right).$$

The desired conclusion follows. \square

³Which does not necessarily requires the convergence of the sequence $\{\|g_k - G_k\|\}$ to zero.

5 Conclusions and perspectives

We have introduced the ADAGB2 algorithm for bound-constrained stochastic optimization which is also capable of using available second-order information. When second-order information is not used and the problem is unconstrained, ADAGB2 subsumes the Adagrad algorithm.

We have shown that, given $\delta \in (0, 1)$ and $\epsilon \in (0, 1]$, the ADAGB2 algorithm needs at most $\mathcal{O}(\epsilon^{-2})$ iterations to ensure an ϵ -approximate first-order critical point of the bound-constrained problem (1) with probability at least $1 - \delta$, provided the average RMSE β_k is sufficiently small. Should this condition fail, we have shown that the optimality default is bounded above for large k by the limit superior of the average oracle’s RMSE.

We have also discussed the relation between the standard optimality measure for bound-constrained induced by the projected gradient in the stochastic case, and have shown that, in general, the unbiased nature of the gradient oracle is not sufficient to ensure convergence on the true constrained problem at the optimal rate, even if such a convergence occurs for the approximate one.

Some interesting theoretical questions remain open at this point, including the extension of the results to (L_0, L_1) -smooth functions, the use of momentum, ensuring second-order optimality and the handling of more general constraints.

Acknowledgements

Serge Gratton and Philippe Toint are grateful to Defeng Sun, Xiaojun Chen and Zaikun Zhang of the Department of Applied Mathematics of the Hong-Kong Polytechnic University for their support during a research visit in the fall 2024. Philippe Toint also acknowledges the support of DIEF (UNIFI) for a visit in October 2024. The research of Stefania Bellavia and Benedetta Morini was partially supported by INDAM-GNCS through Progetti di Ricerca 2023 and by PNRR - Missione 4 Istruzione e Ricerca - Componente C2 Investimento 1.1, Fondo per il Programma Nazionale di Ricerca e Progetti di Rilevante Interesse Nazionale (PRIN) funded by the European Commission under the NextGeneration EU programme, project “Advanced optimization METHods for automated central veIn Sign detection in multiple sclerosis from magneTic resonAnce imaging (AMETISTA)”, code: P2022J9SNP, MUR D.D. financing decree n. 1379 of 1st September 2023 (CUP E53D23017980001), project “Numerical Optimization with Adaptive Accuracy and Applications to Machine Learning”, code: 2022N3ZNAX MUR D.D. financing decree n. 973 of 30th June 2023 (CUP B53D23012670006), and by Partenariato esteso FAIR “Future Artificial Intelligence Research” SPOKE 1 Human-Centered AI. Obiettivo 4, Project “Mathematical and Physical approaches to innovative Machine Learning technologies (MaPLe)”, Codice Identificativo EP_FAIR_002, CUP B93C23001750006.

References

- [1] A. Attia and T. Koren. SGD with AdaGrad stepsizes: full adaptivity with high probability to unknown parameters, unbounded gradients and affine variance. In *International Conference on Machine Learning*, 2023.
- [2] S. Bellavia, B. Morini, and M. Yousefi. Fully stochastic trust-region methods with Barzilai-Borwein steplengths. arXiv:2412.12180, 2024.
- [3] A. Berahas, L. Cao, and K. Scheinberg. Global convergence rate analysis of a generic line search algorithm with noise. *SIAM Journal on Optimization*, 31:1489–1518, 2021.
- [4] A.S. Berahas, M. Jahani, P. Richtárik, and M. Takáč. Quasi-Newton methods for machine learning: forget the past, just sample. *Optimization Methods and Software*, 37(5):1668–1704, 2022.
- [5] J. Blanchet, C. Cartis, M. Menickelly, and K. Scheinberg. Convergence rate analysis of a stochastic trust region method via supermartingales. *INFORMS Journal on Optimization*, 1(2):92–119, 2019.
- [6] L. Bottou, F. E. Curtis, and J. Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311, 2018.
- [7] S. Cai, Z. Mao, Z. Wang, M. Yin, and G. E. Karniadakis. Physics-informed neural networks (PINNs) for fluid mechanics: A review. *Acta Mechanica Sinica*, 37(12):1727–1738, 2021.
- [8] I. Chatzigeorgiou. Bounds on the Lambert function and their application to the outage analysis of user cooperation. *IEEE Communications Letters*, 17(8):1505–1508, 2013.
- [9] G. Ciaramella, F. Nobile, and T. Vanzan. A multigrid solver for PDE-constrained optimization with uncertain inputs. *Journal of Scientific Computing*, 101(1):13, 2024.
- [10] A. R. Conn, N. I. M. Gould, and Ph. L. Toint. *Trust-Region Methods*. Number 1 in MOS-SIAM Optimization Series. SIAM, Philadelphia, USA, 2000.

- [11] R. M. Corless, G. H. Gonnet, D. E. Hare, D. J. Jeffrey, and D. E. Knuth. On the Lambert W function. *Advances in Computational Mathematics*, 5:329–359, 1996.
- [12] F. Curtis, V. Kungurtsev, D. Robinson, and Q. Wang. A stochastic-gradient-based interior-point algorithm for solving smooth bound-constrained optimization problems. arXiv:3204.14907, 2023.
- [13] F.E. Curtis and R. Shi. A fully stochastic second-order trust region method. *Optimization Methods and Software*, 37(3):844–877, 2022.
- [14] J. de Faria, R. Assunção, and F. Murai. Fisher scoring method for neural networks optimization. In *SIAM International Conference on Data Mining (SDM)*, pages 748–756. SIAM, 2023.
- [15] A. Défossez, L. Bottou, F. Bach, and N. Usunier. A simple convergence proof for Adam and Adagrad. *Transactions on Machine Learning Research*, October 2022.
- [16] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(7):2121–2159, 2011.
- [17] M. Faw, L. Rout, C. Caramanis, and S. Shakkottai. Beyond uniform smoothness: a stopped analysis of adaptive SGD. In *Conference on Learning Theory*, 2023.
- [18] M. Faw, I. Tziotis, C. Caramanis, A. Mokhtari, S. Shakkottai, and R. Ward. The power of adaptivity in SGD: self-tuning step sizes with unbounded gradients and affine variance. In *Conference on Learning Theory*, 2022.
- [19] S. Gratton, S. Jerad, and Ph. L. Toint. First-order objective-function-free optimization algorithms and their complexity. arXiv:2203.01757, 2022.
- [20] S. Gratton, S. Jerad, and Ph. L. Toint. Complexity of Adagrad and other first-order methods for nonconvex optimization problems with bounds constraints. arXiv:2406.15793, 2024.
- [21] S. Gratton, S. Jerad, and Ph. L. Toint. Parametric complexity analysis for a class of first-order Adagrad-like algorithms. *Optimization Methods and Software*, (to appear), 2025.
- [22] Y. Hong and J. Lin. Revisiting convergence of Adagrad with relaxed assumptions. arXiv:2403.13794v2, 2024.
- [23] R. Jiang, D. Maladkar, and A. Mokhtari. Convergence analysis of adaptive gradient methods under refined smoothness and noise assumptions. arXiv:2406.04592, 2024.
- [24] A. Kavis, K. Y. Levy, and V. Cevher. High probability bounds for a class of nonconvex algorithms with AdaGrad stepsize. In *International Conference on Learning Representations*, 2022.
- [25] A. Kopaničáková and R. Krause. A recursive multilevel trust region method with application to fully monolithic phase-field models of brittle fracture. *Computer Methods in Applied Mechanics and Engineering*, 360:112720, 2020.
- [26] R. Krause. A Nonsmooth Multiscale Method for Solving Frictional Two-Body Contact Problems in 2D and 3D with Multigrid Efficiency. *SIAM Journal on Scientific Computing*, 31(2):1399–1423, 2009.
- [27] B. Leimkuhler, T. J. Vlaar, T. Pouchon, and A. Storkey. Better training using weight-constrained stochastic dynamics. In *International Conference on Machine Learning*, pages 6200–6211. PMLR, 2021.
- [28] X. Li and F. Orabona. On the convergence of stochastic gradient descent with adaptive stepsizes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, page 983–992, 2019.
- [29] Z. Liu, T. D. Nguyen, A. Ene, and H. Nguyen. On the convergence of Adagrad(Norm) on \mathbb{R}^d : Beyond convexity, non-asymptotic rate and acceleration. In *International Conference on Learning Representations*, 2023.
- [30] L. Lu, R. Pestourie, W. Yao, Z. Wang, F. Verdugo, and S. G. Johnson. Physics-informed neural networks with hard constraints for inverse design. *SIAM Journal on Scientific Computing*, 43(6):B1105–B1132, 2021.
- [31] T. Steihaug. The conjugate gradient method and trust regions in large scale optimization. *SIAM Journal on Numerical Analysis*, 20(3):626–637, 1983.
- [32] C. Tan, S. Ma, Y. H. Dai, and Y. Qian. Barzilai-Borwein step size for stochastic gradient descent. In *30th Conference on Neural Information Processing Systems (NIPS 2016)*, 2016.
- [33] Ph. L. Toint. Towards an efficient sparsity exploiting Newton method for minimization. In I. S. Duff, editor, *Sparse Matrices and Their Uses*, pages 57–88, London, 1981. Academic Press.
- [34] B. Wang, H. Zhang, Z. Ma, and W. Chen. Convergence of AdaGrad for non-convex objectives: simple proofs and relaxed assumptions. In *Conference on Learning Theory*, 2023.
- [35] L. Wang, H. Wu, and I.A. Matveev. Stochastic gradient method with Barzilai-Borwein step for unconstrained nonlinear optimization. *Journal of Computer and Systems Sciences International*, 60(1):75–86, 2021.
- [36] S. Wang, X. Yu, and P. Perdikaris. When and why PINNs fail to train: a neural tangent kernel perspective, 2022.
- [37] R. Ward, X. Wu, and L. Bottou. Adagrad stepsizes: sharp convergence over nonconvex landscapes. *Journal of Machine Learning Research*, 21(1):9047–9076, 2020.
- [38] T. Zhang and Z. Xu. Convergence of stochastic gradient descent for non-convex problems. In *Proceedings for the COLT Conference*, 2012.