

Adaptivity of averaged stochastic gradient descent to local strong convexity for logistic regression

Francis Bach

► To cite this version:

Francis Bach. Adaptivity of averaged stochastic gradient descent to local strong convexity for logistic regression. 2013. hal-00804431v2

HAL Id: hal-00804431 https://hal.science/hal-00804431v2

Preprint submitted on 26 Oct 2013 (v2), last revised 15 Mar 2014 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Adaptivity of averaged stochastic gradient descent to local strong convexity for logistic regression

Francis Bach

FRANCIS.BACH@ENS.FR

INRIA - Sierra Project-team Département d'Informatique de l'Ecole Normale Supérieure Paris, France

Editor:

Abstract

In this paper, we consider supervised learning problems such as logistic regression and study the stochastic gradient method with averaging, in the usual stochastic approximation setting where observations are used only once. We show that after N iterations, with a constant step-size proportional to $1/R^2\sqrt{N}$ where N is the number of observations and R is the maximum norm of the observations, the convergence rate is always of order $O(1/\sqrt{N})$, and improves to $O(R^2/\mu N)$ where μ is the lowest eigenvalue of the Hessian at the global optimum (when this eigenvalue is greater than R^2/\sqrt{N}). Since μ does not need to be known in advance, this shows that averaged stochastic gradient is adaptive to *unknown local* strong convexity of the objective function. Our proof relies on the generalized self-concordance properties of the logistic loss and thus extends to all generalized linear models with uniformly bounded features.

Keywords: Stochastic approximation, logistic regression, self-concordance

1. Introduction

The minimization of an objective function which is only available through unbiased estimates of the function values or its gradients is a key methodological problem in many disciplines. Its analysis has been attacked mainly in three scientific communities: stochastic approximation (Fabian, 1968; Ruppert, 1988; Polyak and Juditsky, 1992; Kushner and Yin, 2003; Broadie et al., 2009), optimization (Nesterov and Vial, 2008; Nemirovski et al., 2009), and machine learning (Bottou and Le Cun, 2005; Shalev-Shwartz et al., 2007; Bottou and Bousquet, 2008; Shalev-Shwartz and Srebro, 2008; Shalev-Shwartz et al., 2009; Duchi and Singer, 2009; Xiao, 2010). The main algorithms which have emerged are stochastic gradient descent (a.k.a. Robbins-Monro algorithm), as well as a simple modification where iterates are averaged (a.k.a. Polyak-Ruppert averaging).

For convex optimization problems, the convergence rates of these algorithms depends primarily on the potential *strong convexity* of the objective function (Nemirovsky and Yudin, 1983). For μ strongly convex functions, after *n* iterations (i.e., *n* observations), the optimal rate of convergence of function values is $O(1/\mu n)$ while for convex functions the optimal rate is $O(1/\sqrt{n})$, both of them achieved by averaged stochastic gradient with step size respectively proportional to $1/\mu n$ or $1/\sqrt{n}$ (Nemirovsky and Yudin, 1983; Agarwal et al., 2012). For smooth functions, averaged stochastic gradient with step sizes proportional to $1/\sqrt{n}$ achieves them up to logarithmic terms (Bach and Moulines, 2011).

Convex optimization problems coming from supervised machine learning are typically of the form $f(\theta) = \mathbb{E}[\ell(y, \langle \theta, x \rangle)]$, where $\ell(y, \langle \theta, x \rangle)$ is the loss between the response $y \in \mathbb{R}$ and the prediction $\langle \theta, x \rangle \in \mathbb{R}$, where x is the input data in a Hilbert space \mathcal{H} and linear predictions parameterized by $\theta \in \mathcal{H}$ are considered. They may or may not have strongly convex objective functions. This most often depends on (a) the correlations between covariates x, and (b) the strong convexity of the loss function ℓ . The logistic loss $\ell : u \mapsto \log(1 + e^{-u})$ is not strongly convex unless restricted to a compact set; moreover, in the sequential observation model, the correlations are not known at training time. Therefore, many theoretical results based on strong convexity do not apply (adding a squared norm $\frac{\mu}{2} ||\theta||^2$ is a possibility, however, in order to avoid adding too much bias, μ has to be small and typically much smaller than $1/\sqrt{n}$, which then makes all strongly-convex bounds vacuous). The goal of this paper is to show that with proper assumptions, namely self-concordance, one can readily obtain favorable theoretical guarantees for logistic regression, namely a rate of the form $O(R^2/\mu n)$ where μ is the lowest eigenvalue of the Hessian at the global optimum, without any exponentially increasing constant factor.

Another goal of this paper is to design an algorithm and provide an analysis that benefit from *hidden* local strong convexity without requiring to know the local strong convexity constant in advance. In smooth situations, the results of Bach and Moulines (2011) imply that the averaged stochastic gradient method with step sizes of the form $O(1/\sqrt{n})$ is adaptive to the strong convexity of the problem. However the dependence in μ in the strongly convex case is of the form $O(1/\mu^2 n)$, which is sub-optimal. Moreover, the final rate is rather complicated, notably because all possible step-sizes are considered. Finally, it does not apply here because even in low-correlation settings, the objective function of logistic regression cannot be globally strongly convex.

In this paper, we provide an analysis for stochastic gradient with averaging for generalized linear models such as logistic regression, with a step size proportional to $1/R^2\sqrt{n}$ where R is the radius of the data and n the number of observations, showing such adaptivity. In particular, we show that the algorithm can adapt to the *local* strong-convexity constant, i.e., the lowest eigenvalue of the Hessian at the optimum. The analysis is done for a finite horizon N and a constant step size decreasing in N as $1/R^2\sqrt{N}$, since the analysis is then slightly easier, though (a) a decaying stepsize could be considered as well, and (b) it could be classically extended to varying step-sizes by a doubling trick (Hazan and Kale, 2001).

2. Stochastic approximation for generalized linear models

In this section, we present the assumptions our work relies on, as well as related work.

2.1 Assumptions

Throughout this paper, we make the following assumptions. We consider a function f defined on a Hilbert space \mathcal{H} , and an increasing family of σ -fields $(\mathcal{F}_n)_{n \ge 1}$; we assume that we are given a

deterministic $\theta_0 \in \mathcal{H}$, and a sequence of functions $f_n : \mathcal{H} \to \mathbb{R}$, for $n \ge 1$. We make the following assumptions, for a certain R > 0:

- (A1) Convexity and differentiability of f: f is convex and three-times differentiable.
- (A2) Generalized self-concordance of f (Bach, 2010): for all $\theta_1, \theta_2 \in \mathcal{H}$, the function $\varphi : t \mapsto f[\theta_1 + t(\theta_2 \theta_1)]$ satisfies: $\forall t \in \mathbb{R}, |\varphi'''(t)| \leq R ||\theta_1 \theta_2||\varphi''(t)$.
- (A3) Attained global minimum: f has a global minimum attained at $\theta_* \in \mathcal{H}$.
- (A4) Lipschitz-continuity of f_n and f: all gradients of f and f_n are bounded by R, that is, for all $\theta \in \mathcal{H}$,

 $||f'(\theta)|| \leq R$ and $\forall n \geq 1$, $||f'_n(\theta)|| \leq R$ almost surely.

- (A5) Adapted measurability: $\forall n \ge 1$, f_n is \mathcal{F}_n -measurable.
- (A6) Unbiased gradients: $\forall n \ge 1$, $\mathbb{E}(f'_n(\theta_{n-1})|\mathcal{F}_{n-1}) = f'(\theta_{n-1})$.
- (A7) Stochastic gradient recursion: $\forall n \ge 1$, $\theta_n = \theta_{n-1} \gamma_n f'_n(\theta_{n-1})$, where $(\gamma_n)_{n\ge 1}$ is a deterministic sequence.

In this paper, we will also consider the averaged iterate $\bar{\theta}_n = \frac{1}{n} \sum_{k=0}^{n-1} \theta_k$, which may be trivially computed on-line through the recursion $\bar{\theta}_n = \frac{1}{n} \theta_{n-1} + \frac{n-1}{n} \bar{\theta}_{n-1}$.

Among the seven assumptions above, the non-standard one is (A2): the notion of self-concordance is an important tool for convex optimization and in particular for the study of Newton's method (Nesterov and Nemirovskii, 1994). It corresponds to having the third derivative bounded by the $\frac{3}{2}$ th power of the second derivative. For machine learning, Bach (2010) has generalized the notion of self-concordance by removing the $\frac{3}{2}$ -th power, so that it is applicable to cost functions arising from probabilistic modeling, as shown below.

Our set of assumptions corresponds to the following examples (with i.i.d. data, and \mathcal{F}_n equal to the σ -field generated by $x_1, y_1, \ldots, x_n, y_n$):

- Logistic regression: $f_n(\theta) = \log(1 + \exp(-y_n \langle x_n, \theta \rangle))$, with data x_n uniformly almost surely bounded by R and $y_n \in \{-1, 1\}$. Note that this includes other binary classification losses, such as $f_n(\theta) = -y_n \langle x_n, \theta \rangle + \sqrt{1 + \langle x_n, \theta \rangle^2}$.
- Generalized linear models with uniformly bounded features: $f_n(\theta) = -\langle \theta, \Phi(x_n, y_n) \rangle + \log \int h(y) \exp (\langle \theta, \Phi(x_n, y) \rangle) dy$, with $\Phi(x_n, y) \in \mathcal{H}$ almost surely bounded in norm by R, for all observations x_n and all potential responses y in a measurable space. This includes multinomial regression and conditional random fields (Lafferty et al., 2001).
- **Robust regression**: we may use $f_n(\theta) = \varphi(y_n \langle x_n, \theta \rangle)$, with $\varphi(t) = \log \cosh t = \log \frac{e^t + e^{-t}}{2}$, with a similar boundedness assumption on x_n .

Running-time complexity. The stochastic gradient descent recursion $\theta_n = \theta_{n-1} - \gamma_n f'_n(\theta_{n-1})$ operates in full generality in the potentially infinite-dimensional Hilbert space \mathcal{H} . There are two practical set-ups where this recursion can be implemented. When \mathcal{H} is finite-dimensional with dimension d, then the complexity of a single iteration is O(d), and thus O(dn) after n iterations. When \mathcal{H} is infinite-dimensional, the recursion can be readily implemented when (a) all functions f_n depend on one-dimensional projections $\langle x_n, \theta \rangle$, i.e., are of the form $f_n(\theta) = \varphi_n(\langle x_n, \theta \rangle)$ for certain

BACH

random functions φ_n (e.g., $\varphi_n(u) = \ell(y_n, u)$ in machine learning), and (b) all scalar products $K_{ij} = \langle x_i, x_j \rangle$ between x_i and x_j , for $i, j \ge 1$, can be computed. This may be done through the classical application of the "kernel trick" (Schölkopf and Smola, 2001; Shawe-Taylor and Cristianini, 2004): if $\theta_0 = 0$, we may represent θ_n as a linear combination of vectors x_1, \ldots, x_n , i.e., $\theta_n = \sum_{i=1}^n \alpha_i x_i$, and the recursion may be written in terms of the weights α_n , through

$$\alpha_n = -\gamma_n x_n \varphi'_n \left(\sum_{i=1}^{n-1} \alpha_i K_{ni} \right).$$

A key element to notice here is that without regularization, the weights α_i corresponding to previous observations remain constant. The overall complexity of the algorithm is $O(n^2)$ times the cost of evaluating a single kernel function. See Bordes et al. (2005); Wang et al. (2012) for approaches aiming at reducing the computational load in this setting. Finally, note that in the kernel setting, the function $f(\theta)$ cannot be strongly convex because the covariance operator of x is typically a compact operator, with a sequence of eigenvalues tending to zero (some regularization is then needed).

2.2 Related work

Non-strongly-convex functions. When only convexity of the objective function is assumed, then several authors (Nesterov and Vial, 2008; Nemirovski et al., 2009; Shalev-Shwartz et al., 2009; Xiao, 2010) have shown that using a step-size proportional to $1/\sqrt{n}$, together with some form of averaging, leads to the minimax optimal rate of $O(1/\sqrt{n})$ (Nemirovsky and Yudin, 1983; Agarwal et al., 2012). Without averaging, the known convergences rates are suboptimal, that is, averaging is key to obtaining the optimal rate (Bach and Moulines, 2011). Note that the smoothness of the loss does not change the rate, but may help to obtain better constants, with the potential use of acceleration (Lan, 2012). Recent work (Bach and Moulines, 2013) has considered algorithms which improve on the rate $O(1/\sqrt{n})$ for smooth self-concordant losses, such as the square and logistic losses. Their analysis relies on some of the results proved in this paper (in particular the high-order bounds in Section 3).

The compactness of the domain is often used within the algorithm (by using orthogonal projections) and within the analysis (in particular to optimize the step size and obtain high-probability bounds). In this paper, we do not make such compactness assumptions, since in a machine learning context, the available bound would be loose and hurt practical performance.

Another difference between several analyses is the use of decaying step sizes of the form $\gamma_n \propto 1/\sqrt{n}$ vs. the use of a constant step size of the form $\gamma \propto 1/\sqrt{N}$ for a finite known horizon N of iterations. The use of a "doubling trick" as done by Hazan and Kale (2001) for strongly convex optimization, where a constant step size is used for iterations between 2^p and 2^{p+1} , with a constant that is proportional to $1/\sqrt{2^p}$, would allow to obtain an anytime algorithm from a finite horizon one. In order to simplify our analysis, we only consider a finite horizon N and a constant step-size that will be proportional to $1/\sqrt{N}$.

Strongly-convex functions. When the function is μ -strongly convex, i.e., $\theta \mapsto f(\theta) - \frac{\mu}{2} ||\theta||^2$ is convex, there are essentially two approaches to obtaining the minimax-optimal rate of $O(1/\mu n)$ (Nemirovsky and Yudin, 1983; Agarwal et al., 2012): (a) using a step size proportional to $1/\mu n$ with

averaging for non-smooth problems (Nesterov and Vial, 2008; Nemirovski et al., 2009; Xiao, 2010; Shalev-Shwartz et al., 2009; Duchi and Singer, 2009; Lacoste-Julien et al., 2012) or a step size proportional to $1/(R^2 + n\mu)$ also with averaging, for smooth problems, where R^2 is the smoothness constant of the loss of a single observation (Le Roux et al., 2012); (b) for smooth problems, using longer step-sizes proportional to $1/n^{\alpha}$ for $\alpha \in (1/2, 1)$ with averaging (Polyak and Juditsky, 1992; Ruppert, 1988; Bach and Moulines, 2011).

Note that the often advocated step size, i.e., of the form C/n where C is larger than $1/\mu$, leads, without averaging to a convergence rate of $O(1/\mu^2 n)$ (Fabian, 1968; Bach and Moulines, 2011), hence with a worse dependence on μ .

The solution (a) requires to have a good estimate of the strong-convexity constant μ , while the second solution (b) does not require to know such estimate and leads to a convergence rate achieving asymptotically the Cramer-Rao lower bound (Polyak and Juditsky, 1992). Thus, this last solution is adaptive to unknown (but positive) amount of strong convexity. However, unless we take the limiting setting $\alpha = 1/2$, it is not adaptive to lack of strong convexity. While the non-asymptotic analysis of Bach and Moulines (2011) already gives a convergence rate in that situation, the bound is rather complicated and also has a suboptimal dependence on μ . Another goal of this paper is to consider a less general result, but more compact and, as already mentioned, a better dependence on the strong convexity constant μ (moreover, as reviewed below, we consider the *local* strong convexity constant, which is much larger).

Finally, note that unless we restrict the support, the objective function for logistic regression cannot be globally strongly convex (since the Hessian tends to zero when $\|\theta\|$ tends to infinity). In this paper we show that stochastic gradient descent with averaging is adaptive to the *local* strong convexity constant, i.e., the lowest eigenvalue of the Hessian of f at the global optimum, without any exponential terms in RD (which would be present if a compact domain of diameter D was imposed and traditional analyses were performed).

Adaptivity to unknown constants. The desirable property of adaptivity to the difficulty of an optimization problem has also been studied in several settings. Gradient descent with constant step size is for example naturally adaptive to the strong convexity of the problem (see, e.g., Nesterov, 2004). In the stochastic context, Juditsky and Nesterov (2010) provide another strategy than averaging with longer step sizes, but for uniform convexity constants.

3. Non-strongly convex analysis

In this section, we study the averaged stochastic gradient method in the non-strongly convex case, i.e., without any (global or local) strong convexity assumptions. We first recall existing results in Section 3.1, that bound the expectation of the excess risk leading to a bound in $O(1/\sqrt{n})$. We then show using martingale moment inequalities how all higher-order moments may be bounded in Section 3.2, still with a rate of $O(1/\sqrt{n})$. However, in Section 3.3, we consider the convergence of the squared gradient, with now a rate of O(1/n). This last result is key to obtaining the adaptivity to local strong convexity in Section 4.

3.1 Existing results

In this section, we review existing results for Lipschitz-continuous non-strongly convex problems (Nesterov and Vial, 2008; Nemirovski et al., 2009; Shalev-Shwartz et al., 2009; Duchi and Singer, 2009; Xiao, 2010). Note that smoothness is not needed here. We consider a constant step size $\gamma_n = \gamma > 0$, for all $n \ge 1$, and we denote by $\bar{\theta}_n = \frac{1}{n} \sum_{k=0}^{n-1} \theta_k$ the averaged iterate.

We prove the following proposition, which provides a bound on the expectation of $f(\bar{\theta}_n) - f(\theta_*)$ that decays at rate $O(\gamma + 1/\gamma n)$, hence the usual choice $\gamma \propto 1/\sqrt{n}$:

Proposition 1 Assume (A1) and (A3-7). With constant step size equal to γ , for any $n \ge 0$, we have:

$$\mathbb{E}f\left(\frac{1}{n}\sum_{k=1}^{n}\theta_{k-1}\right) - f(\theta_*) + \frac{1}{2\gamma n}\mathbb{E}\|\theta_n - \theta_*\|^2 \leqslant \frac{1}{2\gamma n}\|\theta_0 - \theta_*\|^2 + \frac{\gamma}{2}R^2.$$
(1)

Proof We have the following recursion, obtained from the Lipschitz-continuity of f_n :

$$\begin{aligned} \|\theta_n - \theta_*\|^2 &= \|\theta_{n-1} - \theta_*\|^2 - 2\gamma \langle \theta_{n-1} - \theta_*, f'_n(\theta_{n-1}) \rangle + \gamma^2 \|f'_n(\theta_{n-1})\|^2 \\ &\leqslant \|\theta_{n-1} - \theta_*\|^2 - 2\gamma \langle \theta_{n-1} - \theta_*, f'(\theta_{n-1}) \rangle + \gamma^2 R^2 + M_n, \end{aligned}$$

with

$$M_n = -2\gamma \langle \theta_{n-1} - \theta_*, f'_n(\theta_{n-1}) - f'(\theta_{n-1}) \rangle.$$

We thus get, using the classical result from convexity $f(\theta_{n-1}) - f(\theta_*) \leq \langle \theta_{n-1} - \theta_*, f'(\theta_{n-1}) \rangle$:

$$2\gamma [f(\theta_{n-1}) - f(\theta_*)] \leq \|\theta_{n-1} - \theta_*\|^2 - \|\theta_n - \theta_*\|^2 + \gamma^2 R^2 + M_n.$$
(2)

Summing over integers less than n, this implies:

$$\frac{1}{n}\sum_{k=0}^{n-1}f(\theta_k) - f(\theta_*) + \frac{1}{2\gamma n}\|\theta_n - \theta_*\|^2 \leqslant \frac{1}{2\gamma n}\|\theta_0 - \theta_*\|^2 + \frac{\gamma}{2}R^2 + \frac{1}{2\gamma n}\sum_{k=1}^n M_k.$$

We get the desired result by taking expectation in the last inequality, and using the expectation $\mathbb{E}M_k = \mathbb{E}(\mathbb{E}(M_k|\mathcal{F}_{k-1})) = 0$ and $f(\frac{1}{n}\sum_{k=0}^{n-1}\theta_k) \leq \frac{1}{n}\sum_{k=0}^{n-1}f(\theta_k)$.

The following corollary considers a specific choice of the step size (note that the bound is only true for the last iteration):

Corollary 2 Assume (A1) and (A3-7). With constant step size equal to $\gamma = \frac{1}{2R^2\sqrt{N}}$, we have:

$$\forall n \in \{1, \dots, N\}, \ \mathbb{E} \|\theta_n - \theta_*\|^2 \leqslant \|\theta_0 - \theta_*\|^2 + \frac{1}{4R^2},$$
(3)

$$\mathbb{E}f\left(\frac{1}{N}\sum_{k=1}^{N}\theta_{k-1}\right) - f(\theta_*) \leqslant \frac{R^2}{\sqrt{N}}\|\theta_0 - \theta_*\|^2 + \frac{1}{4\sqrt{N}}.$$
(4)

Note that if $\|\theta_0 - \theta_*\|^2$ was known, then a better step-size would be $\gamma = \frac{\|\theta_0 - \theta_*\|}{R\sqrt{N}}$, leading to a convergence rate proportional to $\frac{R\|\theta_0 - \theta_*\|}{\sqrt{N}}$. However, this requires an estimate (simply an upper-bound) of $\|\theta_0 - \theta_*\|^2$, which is typically not available.

We are going to improve this result in several ways:

- All moments of $\|\theta_n \theta_*\|^2$ and $f(\overline{\theta}_n) f(\theta_*)$ will be bounded, leading to a sub-exponential behavior. Note that we do not assume that the iterates are restricted to a predefined bounded set, which is the usual assumption made to derive tail bounds for stochastic approximation (Nesterov and Vial, 2008; Nemirovski et al., 2009; Kakade and Tewari, 2009).
- We are going to show that the squared norm of the gradient at $\bar{\theta}_n = \frac{1}{n} \sum_{k=1}^n \theta_{k-1}$ converges at rate O(1/n), even in the non-strongly convex case. This will allow us to derive finer convergence rates in presence of local strong convexity in Section 4.

3.2 Higher-order bound

In this section, we prove novel higher-order bounds (see the proof in Appendix C, which is based on taking powers of the inequality in Eq. (2) and using martingale moment inequalities), both for any constant step-sizes and then for the specific choice $\gamma = \frac{1}{2R^2\sqrt{N}}$.

Proposition 3 Assume (A1) and (A3-7). With constant step size equal to γ , for any $n \ge 0$ and integer $p \ge 1$, we have:

$$\mathbb{E}\left(2\gamma n \left[f(\bar{\theta}_n) - f(\theta^*)\right] + \|\theta_n - \theta_*\|^2\right)^p \leqslant \left(3\|\theta_0 - \theta_*\|^2 + 20np\gamma^2 R^2\right)^p.$$
(5)

Corollary 4 Assume (A1) and (A3-7). With constant step size equal to $\gamma = \frac{1}{2R^2\sqrt{N}}$, for any integer $p \ge 1$, we have:

$$\forall n \in \{1, \dots, N\}, \ \mathbb{E} \|\theta_N - \theta_*\|^{2p} \leqslant \left[\frac{1}{R_*^2} (3R^2 \|\theta_0 - \theta_*\|^2 + 5p)\right]^p, \tag{6}$$

$$\mathbb{E}\left[f(\bar{\theta}_N) - f(\theta^*)\right]^p \leqslant \left[\frac{1}{\sqrt{N}} \left(3R^2 \|\theta_0 - \theta_*\|^2 + 5p\right)\right]^p.$$
(7)

Having a bound on all moments allows immediately to derive large deviation bounds in the same two cases (by applying Lemma 11 from Appendix A):

Proposition 5 Assume (A1) and (A3-7). With constant step size equal to γ , for any $n \ge 0$ and $t \ge 0$, we have:

$$\mathbb{P}\Big(f(\bar{\theta}_n) - f(\theta_*) \ge 30\gamma R^2 t + \frac{3\|\theta_0 - \theta_*\|^2}{\gamma n}\Big) \le 2\exp(-t),$$
$$\mathbb{P}\Big(\|\theta_n - \theta_*\|^2 \ge 60n\gamma^2 R^2 t + 6\|\theta_0 - \theta_*\|^2\Big) \le 2\exp(-t).$$

Corollary 6 Assume (A1) and (A3-7). With constant step size equal to $\gamma = \frac{1}{2R^2\sqrt{N}}$, for any $t \ge 0$ we have:

$$\mathbb{P}\Big(f(\bar{\theta}_N) - f(\theta_*) \ge \frac{15t}{\sqrt{N}} + \frac{6R^2 \|\theta_0 - \theta_*\|^2}{\sqrt{N}}\Big) \le 2\exp(-t),$$
$$\mathbb{P}\Big(\|\theta_N - \theta_*\|^2 \ge 15R^{-2}t + 6\|\theta_0 - \theta_*\|^2\Big) \le 2\exp(-t).$$

We can make the following observations:

- The results above bounding the norm between the last iterate and a global optimum extends to the averaged iterate.
- The iterates θ_n and $\overline{\theta}_n$ do not necessarily converge to θ_* (note that θ_* may not be unique in general anyway).
- Given that $(\mathbb{E}[f(\bar{\theta}_n) f(\theta_*)]^p)^{1/p}$ is affine in p, we obtain a subexponential behavior, i.e., tail bounds similar to an exponential distribution. The same decay was obtained by Nesterov and Vial (2008) and Nemirovski et al. (2009), but with an extra orthogonal projection step that is equivalent in our setting to know a bound on $\|\theta_*\|$, which is in practise not available.
- The proof of Prop. 3 is rather technical and makes heavy use of martingale moment inequalities.
 A simpler alternative proof has been derived by Bach and Moulines (2013), which uses the Burkholder-Rosenthal-Pinelis inequality (Pinelis, 1994, Theorem 4.1).
- The constants in the bounds of of Prop. 3 (and thus other results as well) could clearly be improved. In particular, we have, for p = 1, 2, 3 (see proof in Appendix E):

$$\mathbb{E}\Big(2\gamma n \big[f(\bar{\theta}_{n}) - f(\theta^{*})\big] + \|\theta_{n} - \theta_{*}\|^{2}\Big) \leq \|\theta_{0} - \theta_{*}\|^{2} + n\gamma^{2}R^{2}, \\
\mathbb{E}\Big(2\gamma n \big[f(\bar{\theta}_{n}) - f(\theta^{*})\big] + \|\theta_{n} - \theta_{*}\|^{2}\Big)^{2} \leq (\|\theta_{0} - \theta_{*}\|^{2} + 9n\gamma^{2}R^{2})^{2}, \\
\mathbb{E}\Big(2\gamma n \big[f(\bar{\theta}_{n}) - f(\theta^{*})\big] + \|\theta_{n} - \theta_{*}\|^{2}\Big)^{3} \leq (\|\theta_{0} - \theta_{*}\|^{2} + 20n\gamma^{2}R^{2})^{3}.$$

3.3 Convergence of gradients

In this section, we prove higher-order bounds on the convergence of the gradient, with an improved rate O(1/n) for $||f'(\bar{\theta}_n)||^2$. In this section, we will need the self-concordance property in Assumption (A2).

Proposition 7 Assume (A1-7). With constant step size equal to γ , for any $n \ge 0$ and integer p, we have:

$$\left(\mathbb{E}\left\|f'\left(\frac{1}{n}\sum_{k=1}^{n}\theta_{k-1}\right)\right\|^{2p}\right)^{1/2p} \leqslant \frac{R}{\sqrt{n}} \left[8\sqrt{p} + \frac{4p}{\sqrt{n}} + 40R^2\gamma p\sqrt{n} + \frac{3}{\gamma\sqrt{n}}\|\theta_0 - \theta_*\|^2 + \frac{3}{\gamma R\sqrt{n}}\|\theta_0 - \theta_*\|\right]$$

$$\tag{8}$$

Corollary 8 Assume (A1-7). With constant step size equal to $\gamma = \frac{1}{2R^2\sqrt{N}}$, for any integer p, we have:

$$\left(\mathbb{E}\left\|f'\left(\frac{1}{N}\sum_{k=1}^{N}\theta_{k-1}\right)\right\|^{2p}\right)^{1/2p} \leqslant \frac{R}{\sqrt{N}} \left[8\sqrt{p} + \frac{4p}{\sqrt{n}} + 20p + 6R^{2}\|\theta_{0} - \theta_{*}\|^{2} + 6R\|\theta_{0} - \theta_{*}\|\right].$$
(9)

We can make the following observations:

- The squared norm of the gradient $||f'(\bar{\theta}_N)||^2$ converges at rate O(1/N).
- Given that $(\mathbb{E} \| f'(\bar{\theta}_N) \|^{2p})^{1/2p}$ is affine in p, we obtain a subexponential behavior for $\| f'(\bar{\theta}_N) \|$, i.e., tail bounds similar to an exponential distribution.

 The proof of Prop. 7 makes use of the self-concordance assumption (that allows to upperbound deviations of gradients by deviations of function values) and of the proof technique of Polyak and Juditsky (1992).

4. Self-concordance analysis for strongly-convex problems

In the previous section, we have shown that $||f'(\bar{\theta}_N)||^2$ is of order O(1/N). If the function f was strongly convex with constant $\mu > 0$, this would immediately lead to the bound $f(\bar{\theta}_N) - f(\theta_*) \leq \frac{1}{2\mu} ||f'(\bar{\theta}_N)||^2$, of order $O(1/\mu N)$. However, because of the Lipschitz-continuity of f on the full Hilbert space \mathcal{H} , it cannot be strongly convex. In this section, we show how the self-concordance assumption may be used to obtain the exact same behavior, but with μ replaced by the *local* strong convexity constant, which is more likely to be strictly positive.

The required property is summarized in the following proposition about (generalized) self-concordant function (see proof in Appendix B.1):

Proposition 9 Let f be a convex three-times differentiable function from \mathcal{H} to \mathbb{R} , such that for all $\theta_1, \theta_2 \in \mathcal{H}$, the function $\varphi : t \mapsto f[\theta_1 + t(\theta_2 - \theta_1)]$ satifies: $\forall t \in \mathbb{R}, |\varphi'''(t)| \leq R ||\theta_1 - \theta_2 ||\varphi''(t)$. Let θ_* be a global minimizer of f and μ the lowest eigenvalue of $f''(\theta_*)$, which is assumed strictly positive.

$$If \frac{\|f'(\theta)\|R}{\mu} \leq \frac{3}{4}, \text{ then } \|\theta - \theta_*\|^2 \leq 4 \frac{\|f'(\theta)\|^2}{\mu^2} \text{ and } f(\theta) - f(\theta_*) \leq 2 \frac{\|f'(\theta)\|^2}{\mu}$$

We may now use this proposition for the averaged stochastic gradient. For simplicity, we only consider the step-size $\gamma = \frac{1}{2R^2\sqrt{N}}$, and the last iterate (see proof in Appendix F):

Proposition 10 Assume (A1-7). Assume $\gamma = \frac{1}{2R^2\sqrt{N}}$. Let $\mu > 0$ be the lowest eigenvalue of the Hessian of f at the unique global optimum θ_* . Then:

$$\mathbb{E}f(\bar{\theta}_N) - f(\theta_*) \leqslant \frac{R^2}{N\mu} \Big(5R\|\theta_0 - \theta_*\| + 15 \Big)^4,$$

$$\mathbb{E}\|\bar{\theta}_N - \theta_*\|^2 \leqslant \frac{R^2}{N\mu^2} \Big(6R\|\theta_0 - \theta_*\| + 21 \Big)^4.$$

We can make the following observations:

- The proof relies on Prop. 9 and requires a control of the probability that $\frac{\|f'(\bar{\theta}_N)\|R}{\mu} \leq \frac{3}{4}$, which is obtained from Prop. 7.
- We conjecture a bound of the form $\left[\frac{R^2}{N\mu}(\Box R \| \theta_0 \theta_* \| + \bigtriangleup \sqrt{p})^4\right]^p$ for the *p*-th order moment of $f(\bar{\theta}_N) f(\theta_*)$, for some scalar constants \Box and \bigtriangleup .
- The new bound now has the term $R \| \theta_0 \theta_* \|$ with a fourth power (compared to the bound in Prop. 1, which has a second power), which typically grows with the dimension of the underlying space (or the slowness of the decay of eigenvalues of the covariance operator when \mathcal{H} is infinite-dimensional). It would be interesting to study whether this dependence can be reduced.

- The key elements in the previous proposition are that (a) the constant μ is the *local* convexity constant, and (b) the step-size does not depend on that constant μ , hence the claimed adaptivity.
- The bounds are only better than the non-strongly-convex bounds from Prop. 1, when the Hessian lowest eigenvalue is large enough, i.e., $\mu R^2 \sqrt{N}$ larger than a fixed constant.
- In the context of logistic regression, even when the covariance matrix of the inputs is invertible, then the only available lower bound on μ is equal to the lowest eigenvalue of the covariance matrix times $\exp(-R\|\theta_*\|)$, which is exponentially small. However, the previous bound is overly pessimistic since it is based on an upper bound on the largest possible value of $\langle x, \theta_* \rangle$. In practice, the actual value of μ is much larger and only a small constant smaller than the lowest eigenvalue of the covariance matrix. In order to assess if this result can be improved, it is interesting to look at the asymptotic result from Polyak and Juditsky (1992) for logistic regression, which leads to a limit rate of 1/n times tr $f''(\theta_*)^{-1}(\mathbb{E}f'_n(\theta_*)f'_n(\theta_*)^{\top})$; note that this rate holds both for the stochastic approximation algorithm and for the global optimum of the training cost, using standard asymptotic statistics results (Van der Vaart, 1998). When the model is well-specified, i.e., the log-odds ratio of the conditional distribution of the label given the input is linear, then $\mathbb{E}f'_n(\theta_*)f'_n(\theta_*)^{\top} = \mathbb{E}f''_n(\theta_*) = f''(\theta_*)$, and the asymptotic rate is exactly d/n, where d is the dimension of \mathcal{H} (which has to be finite-dimensional for the covariance matrix to be invertible). It would be interesting to see if making the extra assumption of well-specification, we can also get an improved *non-asymptotic* result. When the model is mis-specified however, the quantity $\mathbb{E} f'_n(\theta_*) f'_n(\theta_*)^{\top}$ may be large even when $f''(\theta_*)$ is small, and the asymptotic regime does readily lead to an improved bound.

5. Conclusion

In this paper, we have provided a novel analysis of averaged stochastic gradient for logistic regression and related problems. The key aspects of our result are (a) the adaptivity to local strong convexity provided by averaging and (b) the use of self-concordance to obtain a simple bound that does not involve a term which is explicitly exponential in $R \| \theta_0 - \theta_* \|$, which could be obtained by constraining the domain of the iterates.

Our results could be extended in several ways: (a) with a finite and known horizon N, we considered a constant step-size proportional to $1/R^2\sqrt{N}$; it thus seems natural to study the decaying step size $\gamma_n = O(1/R^2\sqrt{n})$, which should, up to logarithmic terms, lead to similar results—and thus likely provide a solution to a a recently posed open problem for online logistic regression (McMahan and Streeter, 2012); (b) an alternative would be to consider a doubling trick where the step-sizes are piecewise constant; also, (c) it may be possible to consider other assumptions, such as exp-concavity (Hazan and Kale, 2001) or uniform convexity (Juditsky and Nesterov, 2010), to derive similar or improved results. Finally, by departing from a plain averaged stochastic gradient recursion, Bach and Moulines (2013) have considered an online Newton algorithm with the same running-time complexity, which leads to a rate of O(1/n) without strong convexity assumptions for logistic regression (though with additional assumptions regarding the distributions of the inputs). It would be interesting to understand if simple assumptions as the ones made in the present paper are possible while preserving the improved convergence rate.

Appendix A. Probability lemmas

In this appendix, we prove simple lemmas relating bounds on moments to tail bounds, with the traditional use of Markov's inequality. See more general results by Boucheron et al. (2013).

Lemma 11 Let X be a non-negative random variable such that for some positive constants A and B, and all $p \in \{1, ..., n\}$,

$$\mathbb{E}X^p \leqslant (A+Bp)^p.$$

Then, if $t \leq \frac{n}{2}$,

$$\mathbb{P}(X \ge 3Bt + 2A) \le 2\exp(-t).$$

Proof We have, by Markov's inequality, for any $p \in \{1, ..., n\}$:

$$\mathbb{P}(X \ge 2Bp + 2A) \leqslant \frac{\mathbb{E}X^p}{(2Bp + 2A)^p} \leqslant \frac{(A + Bp)^p}{(2A + 2Bp)^p} = \exp(-\log(2)p).$$

For $u \in [1, n]$, we consider $p = \lfloor u \rfloor$, so that

$$\mathbb{P}(X \ge 2Bu + 2A) \leqslant \mathbb{P}(X \ge 2Bp + 2A) \leqslant \exp(-\log(2)p) \leqslant 2\exp(-\log(2)u).$$

We take $t = \log(2)u$ and use $2/\log 2 \leq 3$. This is thus valid if $t \leq \frac{n}{2}$.

Lemma 12 Let X be a non-negative random variable such that for some positive constants A, B and C, and for all $p \in \{1, ..., n\}$,

$$\mathbb{E}X^p \leqslant (A\sqrt{p} + Bp + C)^{2p}.$$

Then, if $t \leq n$,

$$\mathbb{P}(X \ge (2A\sqrt{t} + 2Bt + 2C)^2) \le 4\exp(-t).$$

Proof We have, by Markov's inequality, for any $p \in \{1, ..., n\}$:

$$\mathbb{P}(X \ge (2A\sqrt{p} + 2Bp + 2C)^2) \le \frac{\mathbb{E}X^p}{(2A\sqrt{p} + 2Bp + 2C)^{2p}} \le \frac{(A\sqrt{p} + Bp + C)^{2p}}{(2A\sqrt{p} + 2Bp + 2C)^{2p}} \le \exp(-\log(4)p)$$

For $u \in [1, n]$, we consider $p = \lfloor u \rfloor$, so that

$$\mathbb{P}(X \ge (2A\sqrt{u} + 2Bu + 2C)^2) \le \mathbb{P}(X \ge (2A\sqrt{u} + 2Bu + 2C)^2) \le \exp(-\log(2)p) \le 4\exp(-\log(4)u).$$

We take $t = \log(4)u$ and use $\log 4 \ge 1$. This is thus valid if $t \le n$.

Appendix B. Self-concordance properties

In this appendix, we show two lemmas regarding our generalized notion of self-concordance, as well as Prop. 9. For more details, see Bach (2010) and references therein.

The following lemma provide an upper-bound on a one-dimensional self-concordant function at a given point which is based on the gradient at this point and the value and the Hessian at the global minimum. This is key to going in Section 4 from a convergence of gradients to a convergence of function values.

Lemma 13 Let $\varphi : [0,1] \to \mathbb{R}$ a strictly convex three-times differentiable function such that for some S > 0, $\forall t \in [0,1]$, $|\varphi'''(t)| \leq S\varphi''(t)$. Assume $\varphi'(0) = 0$, $\varphi''(0) > 0$. Then:

$$\frac{\varphi'(1)}{\varphi''(0)}S \geqslant 1 - e^{-S} \text{ and } \varphi(1) \leqslant \varphi(0) + \frac{\varphi'(1)^2}{\varphi''(0)}(1+S).$$

Moreover, if $\alpha = \frac{\varphi'(1)S}{\varphi''(0)} < 1$, then $\varphi(1) \leq \varphi(0) + \frac{\varphi'(1)^2}{\varphi''(0)} \frac{1}{\alpha} \log \frac{1}{1-\alpha}$. If in addition $\alpha \leq \frac{3}{4}$, then $\varphi(1) \leq \varphi(0) + 2\frac{\varphi'(1)^2}{\varphi''(0)}$ and $\varphi''(0) \leq 2\varphi'(1)$.

Proof By self-concordance, we obtain that the derivative of $u \mapsto \log \varphi''(u)$ is lower-bounded by -S. By integrating between 0 and $t \in [0, 1]$, we get

$$\log \varphi''(t) - \log \varphi''(0) \ge -St$$
, i.e., $\varphi''(t) \ge \varphi''(0)e^{-St}$,

and by integrating between 0 and 1, we obtain (note that we have assumed $\varphi'(0) = 0$):

$$\varphi'(1) \geqslant \varphi''(0) \frac{1 - e^{-S}}{S}.$$
(10)

We then get (with a first inequality from convexity of φ , and the last inequality from $e^S \ge 1 + S$):

$$\varphi(1) - \varphi(0) \leqslant \varphi'(1) \leqslant \varphi'(1) \frac{\varphi'(1)}{\varphi''(0)} \frac{S}{1 - e^{-S}} = \frac{\varphi'(1)^2}{\varphi''(0)} \left(S + \frac{S}{e^S - 1}\right) \leqslant \frac{\varphi'(1)^2}{\varphi''(0)} (1 + S).$$

Eq. (10) implies that $\alpha \ge 1 - e^{-S}$, which implies, if $\alpha < 1$, $S \le \log \frac{1}{1-\alpha}$. This implies that

$$\varphi(1) - \varphi(0) \leqslant \varphi'(1) \frac{\varphi'(1)}{\varphi''(0)} \frac{S}{1 - e^{-S}} \leqslant \frac{\varphi'(1)^2}{\varphi''(0)} \frac{1}{\alpha} \log \frac{1}{1 - \alpha}$$

using the monotinicity of $S \mapsto \frac{S}{1-e^{-S}}$. Finally the last bounds are a consequence of $\frac{S}{\alpha} \leq \frac{1}{\alpha} \log \frac{1}{1-\alpha} \leq 2$, which is valid for $\alpha \leq \frac{3}{4}$.

The following lemma upper-bounds the remainder in the first-order Taylor expansion of the gradient by the remainder in the first-order Taylor expansion of the function. This is important when function values behave well (i.e., converge to the minimal value) while the iterates may not. **Lemma 14** Let f be a convex three-times differentiable function from \mathcal{H} to \mathbb{R} , such that for all $\theta_1, \theta_2 \in \mathcal{H}$, the function $\varphi : t \mapsto f[\theta_1 + t(\theta_2 - \theta_1)]$ satisfies: $\forall t \in \mathbb{R}, |\varphi'''(t)| \leq R ||\theta_1 - \theta_2||\varphi''(t)$. For any $\theta_1, \theta_2 \in H$, we have:

$$\left\|f'(\theta_1) - f'(\theta_2) - f''(\theta_2)(\theta_2 - \theta_1)\right\| \leqslant R \big[f(\theta_1) - f(\theta_2) - \langle f'(\theta_2), \theta_2 - \theta_1 \rangle \big].$$

Proof For a given $z \in \mathcal{H}$ of unit norm, let $\varphi(t) = \langle z, f'(\theta_2 + t(\theta_1 - \theta_2)) - f'(\theta_2) - tf''(\theta_2)(\theta_2 - \theta_1) \rangle$ and $\psi(t) = R[f(\theta_2 + t(\theta_1 - \theta_2)) - f(\theta_2) - t\langle f'(\theta_2), \theta_2 - \theta_1 \rangle]$. We have $\varphi(0) = \psi(0) = 0$. Moreover, we have the following derivatives:

$$\begin{split} \varphi'(t) &= \langle z, f''(\theta_2 + t(\theta_1 - \theta_2)) - f''(\theta_2), \theta_1 - \theta_2 \rangle \\ \varphi''(t) &= f'''(\theta_2 + t(\theta_1 - \theta_2))[z, \theta_1 - \theta_2, \theta_1 - \theta_2] \\ &\leqslant R \|z\|_2 f''(\theta_2 + t(\theta_1 - \theta_2))[\theta_1 - \theta_2, \theta_1 - \theta_2], \text{ using the Appendix A of Bach (2010),} \\ &= R \langle \theta_2 - \theta_1, f''(\theta_2 + t(\theta_1 - \theta_2))(\theta_1 - \theta_2) \rangle \\ \psi'(t) &= R \langle f'(\theta_2 + t(\theta_1 - \theta_2)) - f'(\theta_2), \theta_1 - \theta_2 \rangle \\ \psi''(t) &= R \langle \theta_2 - \theta_1, f''(\theta_2 + t(\theta_1 - \theta_2))(\theta_1 - \theta_2) \rangle, \end{split}$$

where $f'''(\theta)$ is the third order tensor of third derivatives. This leads to $\varphi'(0) = \psi'(0) = 0$ and $\varphi''(t) \leq \psi''(t)$. We thus have $\varphi(1) \leq \psi(1)$ by integrating twice, which leads to the desired result by maximizing with respect to z.

B.1 Proof of Prop. 9

.. /

We follow the standard proof techniques in self-concordant analysis and define an appropriate function of a single real variable and apply simple lemmas like the ones above.

Define $\varphi : t \mapsto f[\theta_* + t(\theta - \theta_*)] - f(\theta_*)$. We have $\begin{aligned} \varphi'(t) &= \langle f'[\theta_* + t(\theta - \theta_*)], \theta - \theta_* \rangle \\ \varphi''(t) &= \langle \theta - \theta_*, f''[\theta_* + t(\theta - \theta_*)](\theta - \theta_*) \rangle \end{aligned}$

$$\varphi^{\prime\prime\prime}(t) = \langle \theta - \theta_*, f \ [\theta_* + t(\theta - \theta_*)](\theta - \theta_*) \rangle$$

$$\varphi^{\prime\prime\prime}(t) = f^{\prime\prime\prime}[\theta_* + t(\theta - \theta_*)][\theta - \theta_*, \theta - \theta_*, \theta - \theta_*].$$

We thus have: $\varphi(0) = \varphi'(0) = 0, \ 0 \leq \varphi'(1) = \langle f'(\theta), \theta - \theta_* \rangle \leq ||f'(\theta)|| ||\theta - \theta_*||, \ \varphi''(0) = \langle \theta - \theta_*, f''(\theta_*)(\theta - \theta_*) \rangle \geq \mu ||\theta - \theta_*||^2$, and $\varphi(t) \geq 0$ for all $t \in [0, 1]$. Moreover, $\varphi'''(t) \leq R ||\theta - \theta_*|| \varphi''(t)$ for all $t \in [0, 1]$, i.e., Lemma 13 applies with $S = R ||\theta - \theta_*||$. This leads to the desired result, with $\alpha = \frac{\varphi'(1)S}{\varphi''(0)} \leq \frac{||f'(\theta)||R}{\mu}$. Note that we also have (using the second inequality in Lemma 13), for all $\theta \in \mathcal{H}$ (and without any assumption on θ):

$$f(\theta) - f(\theta_*) \leq (1 + R \| \theta - \theta_* \|) \frac{\| f'(\theta) \|^2}{\mu}.$$

Appendix C. Proof of Prop. 3

We consider a direct proof based on taking powers of the inequality in Eq. (2), and then using the appropriate martingale properties. The proof is rather technical and does not use any known martingale inequalities about p-th order moments. A simpler alternative proof has been derived by Bach and Moulines (2013), which uses the Burkholder-Rosenthal-Pinelis inequality (Pinelis, 1994, Theorem 4.1).

The proof works as follows: (a) derive a recursion between the p-th moments and the lower-order moments and (c) prove the result by induction on p. Note that we have to treat separately small values on n in the recursion, which do using almost sure bounds in Appendix C.2.

C.1 Derivation of recursion

From the proof of Prop. 1, we have the recursion:

$$2\gamma [f(\theta_{n-1}) - f(\theta_*)] + \|\theta_n - \theta_*\|^2 \leqslant \|\theta_{n-1} - \theta_*\|^2 + \gamma^2 R^2 + M_n,$$

with

$$M_n = -2\gamma \langle \theta_{n-1} - \theta_*, f'_n(\theta_{n-1}) - f'(\theta_{n-1}) \rangle.$$

This leads to, by summing from 1 to n, and using the convexity of f:

$$2\gamma n f\left(\frac{1}{n}\sum_{k=1}^{n}\theta_{k-1}\right) - 2\gamma n f(\theta^*) + \|\theta_n - \theta_*\|^2 \leqslant A_n,$$

with $A_n = \|\theta_0 - \theta_*\|^2 + n\gamma^2 R^2 + \sum_{k=1}^n M_k \ge 0$, or defined through the recursion $A_n = A_{n-1} + \gamma^2 R^2 + M_n$, with $A_0 = \|\theta_* - \theta_0\|^2$. Note that $\mathbb{E}(M_k | \mathcal{F}_{k-1}) = 0$ and $|M_k| \le 4\gamma R \|\theta_{k-1} - \theta_*\| \le 4\gamma R A_{k-1}^{1/2}$ almost surely. This leads to, by using the binomial expansion formula:

$$\begin{aligned}
A_n^p &\leq \left(A_{n-1} + \gamma^2 R^2 + M_n\right)^p &= \sum_{k=0}^p \binom{p}{k} \left(A_{n-1} + \gamma^2 R^2\right)^{p-k} M_n^k \\
&\leq \left(A_{n-1} + \gamma^2 R^2\right)^p + p \left(A_{n-1} + \gamma^2 R^2\right)^{p-1} M_n + \sum_{k=2}^p \binom{p}{k} \left(A_{n-1} + \gamma^2 R^2\right)^{p-k} \left(4\gamma R A_{n-1}^{1/2}\right)^k.
\end{aligned}$$

This leads to, using $E(M_n | \mathcal{F}_{n-1}) = 0$, upper bounding $\gamma^2 R^2$ by $4\gamma^2 R^2$, and using the binomial expansion formula several times:

$$\mathbb{E}[A_n^p|\mathcal{F}_{n-1}] \leq (A_{n-1} + 4\gamma^2 R^2)^p + \sum_{k=2}^p \binom{p}{k} (A_{n-1} + 4\gamma^2 R^2)^{p-k} (4\gamma R A_{n-1}^{1/2})^k \\
= (A_{n-1} + 4\gamma^2 R^2 + 4\gamma R A_{n-1}^{1/2})^p - 4\gamma R p (A_{n-1} + 4\gamma^2 R^2)^{p-1} A_{n-1}^{1/2} \\
= (A_{n-1}^{1/2} + 2\gamma R)^{2p} - 4\gamma R p (A_{n-1} + 4\gamma^2 R^2)^{p-1} A_{n-1}^{1/2} \\
= \sum_{k=0}^{2p} \binom{2p}{k} A_{n-1}^{k/2} (2\gamma R)^{2p-k} - 4\gamma R p A_{n-1}^{1/2} \sum_{k=0}^{p-1} \binom{p-1}{k} A_{n-1}^k (2\gamma R)^{2(p-1-k)} \\
= \sum_{k=0}^{2p} A_{n-1}^{k/2} (2\gamma R)^{2p-k} C_k,$$

with the constants C_k defined as:

$$C_{2q} = \binom{2p}{2q} \text{ for } q \in \{0, \dots, p\},$$

$$C_{2q+1} = \binom{2p}{2q+1} - 2p\binom{p-1}{q} \text{ for } q \in \{0, \dots, p-1\}.$$

In particular, $C_0 = 1$, $C_{2p} = 1$, $C_1 = 0$ and $C_{2p-1} = {2p \choose 2p-1} - 2p {p-1 \choose p-1} = 0$.

Our goal is now to bounding the values of C_k to obtain Eq. (13) below. This will be done by bounding the odd-indexed element by the even-indexed elements.

We have, for $q \in \{1, ..., p - 2\}$,

$$C_{2q+1}\frac{2q+1}{2p-2q-1} \leqslant \binom{2p}{2q+1}\frac{2q+1}{2p-2q-1} \\ = \frac{(2p)!}{(2q+1)!(2p-2q-1)!}\frac{2q+1}{2p-2q-1} \\ = \frac{(2p)!}{(2q)!(2p-2q)!}\frac{2p-2q}{2p-2q-1} = \binom{2p}{2q}\frac{2p-2q}{2p-2q-1}.$$
(11)

For the end of the interval above in q, i.e., q = p - 2, we obtain $C_{2q+1} \frac{2q+1}{2p-2q-1} \leq C_{2q} \frac{4}{3}$, while for $q \leq p - 3$, we obtain $C_{2q+1} \frac{2q+1}{2p-2q-1} \leq C_{2q} \frac{6}{5}$. Moreover, for $q \in \{1, \ldots, p-2\}$,

$$C_{2q+1} \frac{2p - 2q - 1}{2q + 1} \leqslant \binom{2p}{2q + 1} \frac{2p - 2q - 1}{2q + 1} \\ = \frac{(2p)!}{(2q + 1)!(2p - 2q - 1)!} \frac{2p - 2q - 1}{2q + 1} \\ = \frac{(2p)!}{(2q + 2)!(2p - 2q - 2)!} \frac{2q + 2}{2q + 1} = \binom{2p}{2q + 2} \frac{2q + 2}{2q + 1}.$$

For the end of the interval above in q, i.e., q = 1, we obtain $C_{2q+1} \frac{2p-2q-1}{2q+1} \leq C_{2q+2} \frac{4}{3}$, while for $q \geq 2$, we obtain $C_{2q+1} \frac{2p-2q-1}{2q+1} \leq C_{2q+2} \frac{6}{5}$.

(12)

We have moreover, by using the bound $2\gamma RA_{n-1}^{1/2} \leq \frac{\alpha}{2}(2\gamma R)^2 + \frac{1}{2\alpha}A_{n-1}$ for $\alpha = \frac{2q+1}{2p-2q-1}$:

$$C_{2q+1}A_{n-1}^{q+1/2}(2\gamma R)^{2p-2q-1}$$

$$= C_{2q+1}A_{n-1}^{q}(2\gamma R)^{2p-2q-2}A_{n-1}^{1/2}(2\gamma R)$$

$$\leqslant C_{2q+1}A_{n-1}^{q}(2\gamma R)^{2p-2q-2}\frac{1}{2}\left[\frac{2q+1}{2p-2q-1}(2\gamma R)^{2} + \frac{2p-2q-1}{2q+1}A_{n-1}\right]$$

$$= \frac{1}{2}C_{2q+1}\frac{2p-2q-1}{2q+1}A_{n-1}^{q+1}(2\gamma R)^{2p-2q-2} + \frac{1}{2}C_{2q+1}\frac{2q+1}{2p-2q-1}A_{n-1}^{q}(2\gamma R)^{2p-2q}.$$

By combining the previous inequality with Eq. (11) and Eq. (12), we get that the terms indexed by 2q + 1 are bounded by the terms indexed by 2q + 2 and 2q. All terms with $q \in \{2, \ldots, p-3\}$

BACH

are expanded with constants $\frac{3}{5}$, while for q = 1 and q = p - 2, this is $\frac{2}{3}$. Overall each even term receives a contribution which is less than $\max\{\frac{6}{5}, \frac{3}{5} + \frac{2}{3}, \frac{2}{3}\} = \frac{19}{15}$. This leads to

$$\sum_{q=1}^{p-2} C_{2q+1} A_{n-1}^{q+1/2} (2\gamma R)^{2p-2q-1} \leqslant \frac{19}{15} \sum_{q=0}^{p-1} C_{2q} A_{n-1}^{q} (2\gamma R)^{2p-2q},$$

leading to the recursion that will allow us to derive our result:

$$\mathbb{E}\left[A_{n}^{p}\big|\mathcal{F}_{n-1}\right] \leqslant A_{n-1}^{p} + \frac{34}{15} \sum_{q=0}^{p-1} \binom{2p}{2q} A_{n-1}^{q} (2\gamma R)^{2p-2q}.$$
(13)

C.2 First bound

In this section, we derive an almost sure bound that will be valid for small n. Since $\|\theta_n - \theta_*\| \le \|\theta_{n-1} - \theta_*\| + \gamma R$ almost surely, we have $\|\theta_n - \theta_*\| \le \|\theta_0 - \theta_*\| + n\gamma R$ for all $n \ge 0$. This in turn implies that

$$\begin{aligned} A_n &\leqslant \|\theta_0 - \theta_*\|^2 + n\gamma^2 R^2 + 4\gamma R \sum_{k=1}^n \|\theta_{k-1} - \theta_*\| \text{ using } |M_k| \leqslant 4\gamma R \|\theta_{k-1} - \theta_*\|, \\ &\leqslant \|\theta_0 - \theta_*\|^2 + n\gamma^2 R^2 + 4\gamma R \sum_{k=1}^n \left[\|\theta_0 - \theta_*\| + (k-1)\gamma R \right] \text{ using the inequality above,} \\ &\leqslant \|\theta_0 - \theta_*\|^2 + n\gamma^2 R^2 + 4\gamma n R \|\theta_0 - \theta_*\| + 2\gamma^2 R^2 n^2 \text{ by summing over the first } n-1 \text{ integers,} \\ &\leqslant \|\theta_0 - \theta_*\|^2 + n\gamma^2 R^2 + 2\gamma^2 n^2 R^2 + 2 \|\theta_0 - \theta_*\|^2 + 2\gamma^2 R^2 n^2 \text{ using } ab \leqslant \frac{a^2}{2} + \frac{b^2}{2}, \\ &\leqslant 3\|\theta_0 - \theta_*\|^2 + 5n\gamma^2 R^2 \text{ almost surely.} \end{aligned}$$

C.3 Proof by induction

We now proceed by induction on p. If we assume that $\mathbb{E}A_k^q \leq (3\|\theta_0 - \theta_*\|^2 + kq\gamma^2 R^2 A)^q$ for all q < p, and a certain B (which we will choose to be equal to 20). We first note that if $n \leq 4p$, then from Eq. (14), we have

$$\mathbb{E}A_{n}^{p} \leq (3\|\theta_{0} - \theta_{*}\|^{2} + 5n^{2}\gamma^{2}R^{2})^{p} \\ \leq (3\|\theta_{0} - \theta_{*}\|^{2} + 20np\gamma^{2}R^{2})^{p}.$$

Thus, we only need to consider $n \ge 4p$. We then get from Eq. (13):

$$\mathbb{E} \|\theta_n - \theta_*\|^{2p} \leq \|\theta_0 - \theta_*\|^{2p} + \frac{34}{15} \sum_{k=0}^{n-1} \sum_{q=0}^{p-1} \binom{2p}{2q} \mathbb{E} A_k^q (2\gamma R)^{2p-2q}$$

$$\leq \|\theta_0 - \theta_*\|^{2p} + \frac{34}{15} \sum_{k=0}^{n-1} \sum_{q=0}^{p-1} \binom{2p}{2q} (3\|\theta_0 - \theta_*\|^2 + kq\gamma^2 R^2 B)^q (2\gamma R)^{2p-2q},$$

using the induction hypothesis. We may now sum with respect to k:

$$\begin{split} \mathbb{E} \|\theta_n - \theta_*\|^{2p} &\leqslant \|\theta_0 - \theta_*\|^{2p} + \frac{34}{15} \sum_{q=0}^{p-1} \binom{2p}{2q} (2\gamma R)^{2p-2q} \sum_{k=0}^{n-1} \left(3\|\theta_0 - \theta_*\|^2 + kq\gamma^2 R^2 B\right)^q \\ &\leqslant \|\theta_0 - \theta_*\|^{2p} + \frac{34}{15} \sum_{q=0}^{p-1} \binom{2p}{2q} (2\gamma R)^{2p-2q} \sum_{j=0}^q 3^j \|\theta_0 - \theta_*\|^{2j} \binom{q}{j} (q\gamma^2 R^2 B)^{q-j} \frac{n^{q-j+1}}{q-j+1} \\ &\text{using } \sum_{k=0}^{n-1} k^\alpha \leqslant \frac{n^{\alpha+1}}{\alpha+1} \text{ for any } \alpha > 0, \\ &= \|\theta_0 - \theta_*\|^{2p} + \frac{34}{15} \sum_{j=0}^{p-1} 3^j \|\theta_0 - \theta_*\|^{2j} (4\gamma^2 R^2 n)^{p-j} \sum_{q=j}^{p-1} \binom{2p}{2q} \binom{q}{j} (\frac{qB}{4})^{q-j} \frac{n^{q-p+1}}{q-j+1}, \end{split}$$

by changing the order of summations. We now aim to show that it is less than

$$\left(3\|\theta_0 - \theta_*\|^2 + kp\gamma^2 R^2 B\right)^p = 3^p \|\theta_0 - \theta_*\|^{2p} + \sum_{j=0}^{p-1} 3^j \|\theta_0 - \theta_*\|^{2j} (\gamma^2 R^2 n)^{p-j} (Bp)^{p-j} {p \choose j}.$$

By comparing all terms in $\|\theta_0 - \theta_*\|^{2j}$, this is true as soon as for all $j \in \{0, \dots, p-1\}$,

$$\frac{34}{15} \sum_{q=j}^{p-1} \binom{2p}{2q} \binom{q}{j} (qB/4)^{q-j} \frac{1}{q-j+1} \frac{1}{n^{p-q-1}} \leqslant (Bp/4)^{p-j} \binom{p}{j}$$

$$\Leftrightarrow \frac{34}{15} \sum_{k=0}^{p-1-j} \binom{2p}{2k+2} \binom{p-1-k}{j} ((p-1-k)B/4)^{p-1-k-j} \frac{1}{p-k-j} \frac{1}{n^k} \leqslant (Bp/4)^{p-j} \binom{p}{j},$$

obtained by using the change of variable k = p - 1 - q. This is implied by, using $n \ge 4p$:

$$\frac{136}{15} \sum_{k=0}^{p-1-j} B^{-1-k} p^{-k-p+j} \binom{2p}{2k+2} \frac{\binom{p-1-k}{j}}{\binom{p}{j}} (p-1-k)^{p-1-k-j} \frac{1}{p-k-j} \leqslant 1.$$

By expanding the binomial coefficients and simplifying by p - k - j, this is equivalent to

$$\frac{136}{15} \sum_{k=0}^{p-1-j} B^{-1-k} p^{-k-p+j} {2p \choose 2k+2} \frac{(p-1-k)\cdots(p-k-j+1)}{p\cdots(p-j+1)} (p-1-k)^{p-1-k-j} \le 1.$$

We may now write

$$\frac{(p-1-k)\cdots(p-k-j+1)}{p\cdots(p-j+1)} = \frac{(p-1-k)!}{(p-k-j)!}\frac{(p-j)!}{p!} = \frac{(p-1-k)!}{p!}\frac{(p-j)!}{(p-k-j)!}$$
$$= \frac{(p-j)\cdots(p-k-j+1)}{p\cdots(p-k)},$$

so that we only need to show that

$$\frac{136}{15} \sum_{k=0}^{p-1-j} B^{-1-k} p^{-k-p+j} {2p \choose 2k+2} \frac{(p-j)\cdots(p-k-j+1)}{p\cdots(p-k)} (p-1-k)^{p-1-k-j} \leqslant 1.$$

BACH

We have, by bounding all terms then than p by p:

$$\begin{split} &\frac{136}{15} \sum_{k=0}^{p-1-j} A^{-1-k} p^{-k-p+j} {2p \choose 2k+2} \frac{(p-j)\cdots(p-k-j+1)}{p\cdots(p-k)} (p-1-k)^{p-1-k-j} \\ &\leqslant \frac{136}{15} \sum_{k=0}^{p-1-j} A^{-1-k} p^{-k-p+j} {2p \choose 2k+2} \frac{p^k}{p\cdots(p-k)} p^{p-1-k-j} \\ &= \frac{136}{15} \sum_{k=0}^{p-1-j} A^{-1-k} p^{-k-1} {2p \choose 2k+2} \frac{1}{p\cdots(p-k)} \\ &= \frac{136}{15} \sum_{k=0}^{p-1-j} A^{-1-k} \frac{p^{-k-1}}{(2k+2)!} \frac{2p(2p-1)\cdots(2p-2k-1)}{p\cdots(p-k)} \\ &= \frac{136}{15} \sum_{k=0}^{p-1-j} A^{-1-k} \frac{p^{-2-1}2^{2k+2}}{(2k+2)!} \frac{p(p-1/2)\cdots(p-k-1/2)}{p\cdots(p-k)} \\ &\leqslant \frac{136}{15} \sum_{k=0}^{p-1-j} A^{-1-k} \frac{2^{2k+2}}{(2k+2)!} \text{ by associating all } 2k+2 \text{ terms in ratios which are all less than } 1, \\ &\leqslant \frac{136}{15} \sum_{k=0}^{+\infty} \frac{(2/\sqrt{A})^{2k+2}}{(2k+2)!} = \frac{136}{15} \left[\cosh(2/\sqrt{A}) - 1 \right] < 1 \text{ if } A \leqslant 20. \end{split}$$

We thus get the desired result $\mathbb{E}A_n^p \leq (3\|\theta_0 - \theta_*\|^2 + 20np\gamma^2 R^2)^p$, and the proposition is proved by induction.

Appendix D. Proof of Prop. 7

The proof is organized in two parts: first show a bound on the averaged gradient $\frac{1}{n} \sum_{k=1}^{n} f'(\theta_{k-1})$, then relate it to the gradient at the averaged iterate, i.e., $f'(\frac{1}{n} \sum_{k=1}^{n} \theta_{k-1})$, using self-concordance.

D.1 Bound on $\frac{1}{n} \sum_{k=1}^{n} f'(\theta_{k-1})$

We have, following Polyak and Juditsky (1992); Bach and Moulines (2011):

$$f'_n(\theta_{n-1}) = \frac{1}{\gamma}(\theta_{n-1} - \theta_n)$$

which implies, by summing over all integers between 1 and n:

$$\frac{1}{n}\sum_{k=1}^{n}f'(\theta_{k-1}) = \frac{1}{n}\sum_{k=1}^{n}\left[f'(\theta_{k-1}) - f'_{k}(\theta_{k-1})\right] + \frac{1}{\gamma n}(\theta_{0} - \theta_{*}) + \frac{1}{\gamma n}(\theta_{*} - \theta_{n}).$$

We denote $X_k = \frac{1}{n} \left[f'(\theta_{k-1}) - f'_k(\theta_{k-1}) \right] \in \mathcal{H}$. We have: $||X_k|| \leq \frac{2R}{n}$ almost surely and $\mathbb{E}(X_k | \mathcal{F}_{k-1}) = 0$, with $\left(\sum_{k=1}^n \mathbb{E}(||X_k||^2 | \mathcal{F}_{k-1}) \right)^{1/2} \leq \frac{2R}{\sqrt{n}}$. We may thus apply the Burkholder-

Rosenthal-Pinelis inequality (Pinelis, 1994, Theorem 4.1), and get:

$$\left[\mathbb{E}\left\|\frac{1}{n}\sum_{k=1}^{n}\left[f'(\theta_{k-1}) - f'_{k}(\theta_{k-1})\right]\right\|^{2p}\right]^{1/2p} \leq 2p\frac{2R}{n} + \sqrt{2p}\frac{2R}{n^{1/2}}$$

This leads to, using Prop. 3 and Minkowski's inequality:

$$\begin{bmatrix} \mathbb{E} \left\| \frac{1}{n} \sum_{k=1}^{n} f'(\theta_{k-1}) \right\|^{2p} \end{bmatrix}^{1/2p} \leqslant \left[\mathbb{E} \left\| \frac{1}{n} \sum_{k=1}^{n} \left[f'(\theta_{k-1}) - f'_{k}(\theta_{k-1}) \right] \right\|^{2p} \right]^{1/2p} \\ + \frac{1}{\gamma n} \left\| \theta_{0} - \theta_{*} \right\| + \frac{1}{\gamma n} \left[\mathbb{E} \left\| \theta_{*} - \theta_{n} \right\|^{2p} \right]^{1/2p} \\ \leqslant 2p \frac{2R}{n} + \sqrt{2p} \frac{2R}{n^{1/2}} + \frac{1}{\gamma n} \left\| \theta_{0} - \theta_{*} \right\| + \left[\frac{1}{\gamma n} \sqrt{3} \left\| \theta_{0} - \theta_{*} \right\|^{2} + 20np\gamma^{2}R^{2} \right] \\ \leqslant 2p \frac{2R}{n} + \sqrt{2p} \frac{2R}{n^{1/2}} + \frac{1}{\gamma n} \left\| \theta_{0} - \theta_{*} \right\| + \left[\frac{\sqrt{3}}{\gamma n} \left\| \theta_{0} - \theta_{*} \right\| + \frac{1}{\gamma n} \sqrt{20np\gamma}R \right] \\ \leqslant \frac{4pR}{n} + \sqrt{2p} \frac{2R}{n^{1/2}} + \frac{2}{\gamma n} \left\| \theta_{0} - \theta_{*} \right\| + \frac{1}{\gamma n} \sqrt{20np\gamma}R \\ \leqslant \frac{4pR}{n} + \sqrt{p} \frac{R}{\sqrt{n}} \left[2\sqrt{2} + \sqrt{20} \right] + \frac{1 + \sqrt{3}}{\gamma n} \left\| \theta_{0} - \theta_{*} \right\| \\ \leqslant \frac{4pR}{n} + 8\sqrt{p} \frac{R}{\sqrt{n}} + \frac{3}{\gamma n} \left\| \theta_{0} - \theta_{*} \right\|.$$
(15)

D.2 Using self-concordance

Using the self-concordance property of Lemma 14 several times, we obtain:

$$\begin{split} & \left\| \frac{1}{n} \sum_{k=1}^{n} f'(\theta_{k-1}) - f'\left(\frac{1}{n} \sum_{k=1}^{n} \theta_{k-1}\right) \right\| \\ &= \left\| \frac{1}{n} \sum_{k=1}^{n} \left[f'(\theta_{k-1}) - f'(\theta_{*}) - f''(\theta_{*})(\theta_{k-1} - \theta_{*}) \right] - f'\left(\frac{1}{n} \sum_{k=1}^{n} \theta_{k-1}\right) + f'(\theta_{*}) + f''(\theta_{*})\left(\frac{1}{n} \sum_{k=1}^{n} \theta_{k-1} - \theta_{*}\right) \right\| \\ &\leqslant \frac{R}{n} \sum_{k=1}^{n} \left[f(\theta_{k-1}) - f(\theta_{*}) - \langle f'(\theta_{*}), \theta_{k-1} - \theta_{*} \rangle \right] + R \left[f\left(\frac{1}{n} \sum_{k=1}^{n} \theta_{k-1}\right) - f(\theta_{*}) + \left\langle f'(\theta_{*}), \frac{1}{n} \sum_{k=1}^{n} \theta_{k-1} - \theta_{*} \right\rangle \right] \\ &\leqslant 2R \left(\frac{1}{n} \sum_{k=1}^{n} f(\theta_{k-1}) - f(\theta_{*})\right) \text{ using the convexity of } f. \end{split}$$

This leads to, using Prop. 3:

$$\left(\mathbb{E}\left\|\frac{1}{n}\sum_{k=1}^{n}f'(\theta_{k-1}) - f'\left(\frac{1}{n}\sum_{k=1}^{n}\theta_{k-1}\right)\right\|^{2p}\right)^{1/2p} \leq 2R\left(\mathbb{E}\left[\frac{1}{n}\sum_{k=1}^{n}f(\theta_{k-1}) - f(\theta_{*})\right]^{2p}\right)^{1/2p} \leq \frac{2R}{2\gamma n}\left(3\|\theta_{0} - \theta_{*}\|^{2} + 40np\gamma^{2}R^{2}\right).$$
 (16)

Summing Eq. (15) and Eq. (16) leads to the desired result.

Appendix E. Results for small *p*

In Prop. 3, we may replace the bound $3\|\theta_0 - \theta_*\|^2 + 20np\gamma^2 R^2$ with a bound with smaller constants for p = 1, 2, 3 (to be used in proofs of results in Section 4). This is done using the same proof principle but finer derivations, as follows. We denote $\gamma^2 R^2 = b$ and $\|\theta - \theta_*\|^2 = a$, and consider the following inequalities which we have considered in the proof of Prop. 3:

$$\begin{aligned} A_n^p &\leqslant (A_{n-1} + b + M_n)^p \\ M_n &\leqslant 4b^{1/2}A_{n-1}^{1/2} \text{ and } \mathbb{E}(M_n|\mathcal{F}_{n-1}) = 0, \\ A_0 &= a. \end{aligned}$$

We simply take expansions of the *p*-th power above, and sum for all first integers. We have:

$$\begin{split} \mathbb{E}A_n &\leqslant \quad \mathbb{E}A_{n-1} + b \leqslant a + nb, \\ \mathbb{E}A_n^2 &\leqslant \quad \mathbb{E}(A_{n-1}^2 + b^2 + 2bA_{n-1} + M_n^2) \leqslant \mathbb{E}A_{n-1}^2 + 2\mathbb{E}A_{n-1}b + b^2 + 16b\mathbb{E}A_{n-1} \\ &\leqslant \quad a^2 + 18b \left[\sum_{k=0}^{n-1} a + kb\right] + b^2n \leqslant a^2 + 18b[na + \frac{n^2}{2}b] + b^2n \text{ using the result about } \mathbb{E}A_{n-1}, \\ &= \quad a^2 + 18bna + b^2(n+9n^2) \\ &\leqslant \quad (a+9nb)^2. \end{split}$$

We may now pursue for the third order moments:

$$\begin{split} \mathbb{E}A_n^3 &\leqslant \quad \mathbb{E}(A_{n-1}+b)^3 + 3\mathbb{E}(A_{n-1}+b)^2 M_n^2 + 3\mathbb{E}(A_{n-1}+b)^3 M_n + \mathbb{E}M_{n-1}^3 \\ &\leqslant \quad \mathbb{E}(A_{n-1}+b)^3 + 3\mathbb{E}(A_{n-1}+b)^2 16b A_{n-1} + 0 + 64b^{3/2} \mathbb{E}A_{n-1}^{3/2} \\ &\leqslant \quad \mathbb{E}(A_{n-1}^3 + 3\mathbb{E}A_{n-1}^2 b + 3\mathbb{E}A_{n-1}b^2 + b^3) + 3(\mathbb{E}A_{n-1} + b) 16b A_{n-1} + 64b^{3/2} \mathbb{E}A_{n-1}^{3/2} \\ &= \quad (\mathbb{E}A_{n-1}^3 + 3\mathbb{E}A_{n-1}^2 b + 3A_{n-1}b^2 + b^3) + 3(\mathbb{E}A_{n-1} + b) 16b \mathbb{E}A_{n-1} + 32b \mathbb{E}A_{n-1} [2b^{/2}A_{n-1}^{1/2}] \\ &\leqslant \quad (\mathbb{E}A_{n-1}^3 + 3\mathbb{E}A_{n-1}^2 b + 3\mathbb{E}A_{n-1}b^2 + b^3) + 3(\mathbb{E}A_{n-1} + b) 16b \mathbb{E}A_{n-1} + 32\mathbb{E}bA_{n-1} [\frac{A_{n-1}}{4} + 4b] \\ &= \quad \mathbb{E}A_{n-1}^3 + 3\mathbb{E}A_{n-1}^2 b + 3\mathbb{E}A_{n-1}b^2 + b^3) + 3(\mathbb{E}A_{n-1} + b) 16b \mathbb{E}A_{n-1} + 32\mathbb{E}bA_{n-1} [\frac{A_{n-1}}{4} + 4b] \\ &= \quad \mathbb{E}A_{n-1}^3 + \mathbb{E}A_{n-1}^2 b [3 + 48 + 8] + \mathbb{E}A_{n-1}b^2 [3 + 48 + 128] + b^3 \\ &= \quad \mathbb{E}A_{n-1}^3 + 59\mathbb{E}A_{n-1}^2 b + 179\mathbb{E}A_{n-1}b^2 + b^3 \\ &\leqslant \quad a^3 + 59b \Big[\sum_{k=1}^{n-1} a^2 + 18bka + b^2(k + 9k^2)\Big] + 179b^2 \Big[\sum_{k=1}^{n-1} a + kb\Big] + nb^3 \\ &\leqslant \quad a^3 + 59b [na^2 + 9bn^2a + b^2(n^2/2 + 3n^3)] + 179b^2[na + bn^2/2] + nb^3 \\ &= \quad a^3 + 59nba^2 + b^2a [59 \cdot 9n^2 + 179n] + b^3 [59/2 \cdot n^2 + 3 \cdot 59n^3 + 179/2 \cdot n^2 + n] \\ &= \quad a^3 + 59nba^2 + b^2a [531n^2 + 179n] + b^3 [119n^2 + 177n^3 + n] \\ &\leqslant \quad (a + 20nb)^3. \end{split}$$

We then obtain:

$$\mathbb{E}\left[2\gamma n \left[f(\bar{\theta}_{n}) - f(\theta^{*})\right] + \|\theta_{n} - \theta_{*}\|^{2}\right]^{2} \leq \left(\|\theta_{0} - \theta_{*}\|^{2} + 9n\gamma^{2}R^{2}\right)^{2}$$
$$\mathbb{E}\left[2\gamma n \left[f(\bar{\theta}_{n}) - f(\theta^{*})\right] + \|\theta_{n} - \theta_{*}\|^{2}\right]^{3} \leq \left(\|\theta_{0} - \theta_{*}\|^{2} + 20n\gamma^{2}R^{2}\right)^{3}.$$

Appendix F. Proof of Prop. 10

The proof follows from applying self-concordance properties (Prop. 9) applied to $\bar{\theta}_n$. We thus need to provide a control on the probability that $||f'(\bar{\theta}_n)|| \ge \frac{3\mu}{4R}$.

F.1 Tail bound for $||f'(\bar{\theta}_n)||$

We derive a large deviation bound, as a consequence of the bound on all moments of $||f'(\bar{\theta}_n)||$ (Prop. 7) and Lemma 12, that allows to go from moments to tail bounds:

$$\mathbb{P}\bigg(\left\|f'(\bar{\theta}_n)\right\| \ge \frac{2R}{\sqrt{n}} \bigg[10\sqrt{t} + 40R^2\gamma t\sqrt{n} + \frac{3}{\gamma\sqrt{n}}\|\theta_0 - \theta_*\|^2 + \frac{3}{\gamma R\sqrt{n}}\|\theta_0 - \theta_*\|\bigg]\bigg) \le 4\exp(-t).$$

In order to derive the bound above, we need to assume that $p \leq n/4$ (so that $4p/n \leq 2\sqrt{p}/\sqrt{n}$), and thus, when applying Lemma 12, the bound above is valid as long as $t \leq n/4$. It is however valid for all t, because the gradients are bounded by R, and for t > n, we have $\frac{2R}{\sqrt{n}} 10\sqrt{t} \geq R$, and the inequality is satisfied with zero probability.

F.2 Bounding the function values

From Prop. 9, if $||f'(\bar{\theta}_n)|| \ge \frac{3\mu}{4R}$, then $f(\bar{\theta}_n) - f(\theta_*) \le 2\frac{||f'(\bar{\theta}_n)||^2}{\mu}$. This will allow us to derive a tail bound for $f(\bar{\theta}_n) - f(\theta_*)$, for sufficiently small deviations. For larger deviations, we will use the tail bound which does not use strong convexity (Prop. 5).

We consider the event

$$A_t = \left\{ \left\| f'(\bar{\theta}_n) \right\| \leq \frac{2R}{\sqrt{n}} \left[10\sqrt{t} + 40R^2\gamma t\sqrt{n} + \frac{3}{\gamma\sqrt{n}} \|\theta_0 - \theta_*\|^2 + \frac{3}{\gamma R\sqrt{n}} \|\theta_0 - \theta_*\| \right] \right\}.$$

We make the following two assumptions regarding γ and t:

$$10\sqrt{t} + 40R^{2}\gamma t\sqrt{n} \leqslant \frac{2}{3}\frac{3\mu}{4R}\frac{\sqrt{n}}{2R} = \frac{\mu\sqrt{n}}{4R^{2}}$$
(17)
and $\frac{3}{\gamma\sqrt{n}} \|\theta_{0} - \theta_{*}\|^{2} + \frac{3}{\gamma R\sqrt{n}} \|\theta_{0} - \theta_{*}\| \leqslant \frac{1}{3}\frac{3\mu}{4R}\frac{\sqrt{n}}{2R} = \frac{\mu\sqrt{n}}{8R^{2}},$

so that the upper-bound on $||f'(\bar{\theta}_n)||$ in the definition of A_t is less than $\frac{3\mu}{4R}$ (so that we can apply Prop. 9). We thus have:

$$\begin{split} A_t &\subset \left\{ f(\bar{\theta}_n) - f(\theta_*) \leqslant \frac{8R^2}{\mu n} \bigg[10\sqrt{t} + 40R^2\gamma t\sqrt{n} + \frac{3}{\gamma\sqrt{n}} \|\theta_0 - \theta_*\|^2 + \frac{2}{\gamma R\sqrt{n}} \|\theta_0 - \theta_*\| \bigg]^2 \right\} \\ &\subset \left\{ f(\bar{\theta}_n) - f(\theta_*) \leqslant \frac{8R^2}{\mu n} \bigg[10\sqrt{t} + 20\Box t + \Delta \bigg]^2 \bigg\}, \\ \text{with } \Box &= 2\gamma R^2\sqrt{n} \text{ and } \Delta = \frac{3}{\gamma\sqrt{n}} \|\theta_0 - \theta_*\|^2 + \frac{3}{\gamma R\sqrt{n}} \|\theta_0 - \theta_*\|. \end{split}$$

This implies that for all $t \ge 0$, such that $10\sqrt{t} + 20\Box t \le \frac{\mu\sqrt{n}}{4R^2}$, i.e., our assumption in Eq. (17), we may apply the tail bound from Appendix F.1 to get:

$$\mathbb{P}\left(f(\bar{\theta}_n) - f(\theta_*) \ge \frac{8R^2}{\mu n} \left[10\sqrt{t} + 20\Box t + \Delta\right]^2\right) \le 4e^{-t}.$$
(18)

Moreover, we have for all $v \ge 0$ (from Prop. 5):

$$\mathbb{P}\left(f(\bar{\theta}_n) - f(\theta_*) \ge 30\gamma R^2 v + \frac{3\|\theta_0 - \theta_*\|^2}{\gamma n}\right) \le 2\exp(-v).$$
(19)

We may now use the last two inequalities to bound the expectation $\mathbb{E}[f(\bar{\theta}_n) - f(\theta_*)]$. We first express the expectation as an integral of the tail bound and split it into three parts:

$$\mathbb{E}\left[f(\bar{\theta}_{n}) - f(\theta_{*})\right] = \int_{0}^{+\infty} \mathbb{P}\left[f(\bar{\theta}_{n}) - f(\theta_{*}) \ge u\right] du$$

$$= \int_{0}^{\Delta^{2} \frac{8R^{2}}{\mu n}} \mathbb{P}\left[f(\bar{\theta}_{n}) - f(\theta_{*}) \ge u\right] du$$

$$+ \int_{\Delta^{2} \frac{8R^{2}}{\mu n}}^{\frac{8R^{2}}{\mu n} \left(\frac{\mu\sqrt{n}}{4R^{2}} + \Delta\right)^{2}} \mathbb{P}\left[f(\bar{\theta}_{n}) - f(\theta_{*}) \ge u\right] du$$

$$+ \int_{\frac{8R^{2}}{\mu n} \left(\frac{\mu\sqrt{n}}{4R^{2}} + \Delta\right)^{2}}^{+\infty} \mathbb{P}\left[f(\bar{\theta}_{n}) - f(\theta_{*}) \ge u\right] du.$$
(20)

We may now bound the three terms separately. For the first integral, we bound the probability by one to get $\int_{0}^{\Delta^{2}\frac{8R^{2}}{\mu n}} \mathbb{P}[f(\bar{\theta}_{n}) - f(\theta_{*}) \ge u] du \leqslant \Delta^{2}\frac{8R^{2}}{n\mu}.$

For the third term in Eq. (20), we use the tail bound in Eq. (19) to get

$$\int_{\frac{8R^{2}}{\mu n}\left(\frac{\mu\sqrt{n}}{4R^{2}}+\Delta\right)^{2}}^{+\infty} \mathbb{P}\left[f(\bar{\theta}_{n})-f(\theta_{*}) \ge u\right] du$$

$$= \int_{\frac{8R^{2}}{\mu n}\left(\frac{\mu\sqrt{n}}{4R^{2}}+\Delta\right)^{2}-\frac{3}{\gamma n}\|\theta_{0}-\theta_{*}\|^{2}}^{+\infty} \mathbb{P}\left[f(\bar{\theta}_{n})-f(\theta_{*}) \ge u+\frac{3}{\gamma n}\|\theta_{0}-\theta_{*}\|^{2}\right] du$$

$$\leqslant 2\int_{\frac{8R^{2}}{\mu n}\left(\frac{\mu\sqrt{n}}{4R^{2}}+\Delta\right)^{2}-\frac{3}{\gamma n}\|\theta_{0}-\theta_{*}\|^{2}}^{+\infty} \exp\left(-\frac{u}{30\gamma R^{2}}\right) du.$$

We may apply Eq. (19) because

$$\frac{8R^2}{\mu n} \left(\frac{\mu\sqrt{n}}{4R^2} + \Delta\right)^2 - \frac{3}{\gamma n} \|\theta_0 - \theta_*\|^2 \ge \frac{8R^2}{\mu n} \left(\frac{\mu\sqrt{n}}{4R^2} + \Delta\right)^2 - \frac{\mu}{8R^2} \ge \frac{8R^2}{\mu n} \left(\frac{\mu\sqrt{n}}{4R^2}\right)^2 - \frac{\mu}{8R^2} = \frac{3\mu}{8R^2} \ge 0.$$

We can now compute the bound explicitly to get

$$\begin{split} &\int_{\frac{8R^2}{\mu n}\left(\frac{\mu\sqrt{n}}{4R^2}+\Delta\right)^2}^{+\infty} \mathbb{P}\big[f(\bar{\theta}_n)-f(\theta_*) \geqslant u\big] du \\ \leqslant & 60\gamma R^2 \exp\left(-\frac{1}{30\gamma R^2} \bigg[\frac{8R^2}{\mu n} \big(\frac{\mu\sqrt{n}}{4R^2}+\Delta\big)^2 - \frac{3}{\gamma n} \|\theta_0-\theta_*\|^2\bigg]\right) \leqslant 60\gamma R^2 \exp\left(-\frac{1}{30\gamma R^2} \frac{3\mu}{8R^2}\right) \\ \leqslant & 60\gamma R^2 \exp\left(-\frac{\mu}{80\gamma R^4}\right) \leqslant 60\gamma R^2 \frac{80\gamma R^4}{2\mu} \text{ using } e^{-\alpha} \leqslant \frac{1}{2\alpha} \text{ for all } \alpha > 0 \\ = & \frac{2400\gamma^2 R^6}{\mu}. \end{split}$$

We now consider the second term in Eq. (20) for which we will use Eq. (18). We consider the change of variable $u = \frac{8R^2}{\mu n} \left[10\sqrt{t} + 20\Box t + \Delta \right]^2$, for which $u \in \left[\Delta^2 \frac{8R^2}{\mu n}, \frac{8R^2}{\mu n} \left(\frac{\mu\sqrt{n}}{4R^2} + \Delta \right)^2 \right]$ implies $t \in [0, +\infty)$. This implies that

$$\begin{split} &\int_{\Delta^{2} \frac{8R^{2}}{\mu n}}^{\frac{8R^{2}}{\mu n} \left(\frac{\mu\sqrt{n}}{4R^{2}} + \Delta\right)^{2}} \mathbb{P} \left[f(\bar{\theta}_{n}) - f(\theta_{*}) \geqslant u \right] du \\ \leqslant & \int_{0}^{\infty} 4e^{-t} d \left(\frac{8R^{2}}{\mu n} \left[10\sqrt{t} + 20\Box t + \Delta \right]^{2} \right) \\ &= & \frac{32R^{2}}{\mu n} \int_{0}^{\infty} e^{-t} \left(100 + 400\Box^{2}2t + 400\Box\frac{3}{2}t^{1/2} + 20\Delta\frac{1}{2}t^{-1/2} + 40\Delta\Box \right) dt \\ &= & \frac{32R^{2}}{\mu n} \left(100\Gamma(1) + 400\Box^{2}2\Gamma(2) + 400\Box\frac{3}{2}\Gamma(3/2) + 20\Delta\frac{1}{2}\Gamma(1/2) + 40\Delta\Box\Gamma(1) \right) \\ & \text{ with } \Gamma \text{ denoting the Gamma function,} \\ &= & \frac{32R^{2}}{\mu n} \left(100 + 400\Box^{2}2 + 400\Box\frac{3}{2}\frac{1}{2}\sqrt{\pi} + 20\Delta\frac{1}{2}\sqrt{\pi} + 40\Delta\Box \right). \end{split}$$

We may now combine the three bounds to get, from Eq. (20),

$$\mathbb{E} \left[f(\bar{\theta}_n) - f(\theta_*) \right] \leq \Delta^2 \frac{8R^2}{n\mu} + \frac{2400\gamma^2 R^6}{\mu} \\ + \frac{32R^2}{\mu n} \left(100 + 400\Box^2 2 + 400\Box \frac{3}{2} \frac{1}{2}\sqrt{\pi} + 20\Delta \frac{1}{2}\sqrt{\pi} + 40\Delta \Box \right) \\ \leq \frac{32R^2}{n\mu} \left[\frac{\Delta^2}{4} + 75\gamma^2 R^4 n + 100 + 800\Box^2 + 300\Box\sqrt{\pi} + 10\Delta\sqrt{\pi} + 40\Delta \Box \right] .$$

For
$$\gamma = \frac{1}{2R^2\sqrt{N}}$$
, with $\alpha = R \|\theta_0 - \theta_*\|$, $\Box = 1$ and $\Delta = 6\alpha^2 + 6\alpha$, we obtain

$$\mathbb{E}[f(\bar{\theta}_N) - f(\theta_*)] \leq \frac{32R^2}{N\mu} \Big[\frac{1}{4} \Delta^2 + 1451 + 58\Delta \Big]$$

$$\leq \frac{32R^2}{N\mu} \Big[9\alpha^4 + 18\alpha^3 + 9\alpha^2 + 1451 + 348\alpha^2 + 348\alpha \Big]$$

$$\leq \frac{R^2}{N\mu} (625\alpha^4 + 7500\alpha^3 + 33750\alpha^2 + 67500\alpha + 50625) = \frac{R^2}{N\mu} (5\alpha + 15)^4.$$

Note that the previous bound is only valid if $\frac{3}{\gamma\sqrt{n}} \|\theta_0 - \theta_*\|^2 + \frac{3}{\gamma R\sqrt{n}} \|\theta_0 - \theta_*\| \leq \frac{\mu\sqrt{n}}{8R^2}$, i.e., under the condition $6R^2 \|\theta_0 - \theta_*\|^2 + 6R \|\theta_0 - \theta_*\| \leq \frac{\mu\sqrt{N}}{8R^2}$. If the condition is not satisfied, then the bound is still valid because of Prop. 1. We thus obtain the desired result.

F.3 Bound on iterates

Following the same principle as for function values in Appendix F.2, we consider the same event A_t . With the same condition on γ and t, we have:

$$A_t \subset \left\{ \|\bar{\theta}_n - \theta_*\|^2 \leqslant \frac{16R^2}{\mu^2 n} \left[10\sqrt{t} + 20\Box t + \Delta \right]^2 \right\},$$

which leads to the tail bound:

$$\mathbb{P}\left(\|\bar{\theta}_n - \theta_*\|^2 \ge \frac{16R^2}{\mu^2 n} \left[10\sqrt{t} + 20\Box t + \Delta\right]^2\right) \le 4e^{-t}.$$
(21)

We may now split the expectation in three integrals:

$$\mathbb{E}\|\bar{\theta}_n - \theta_*\|^2 = \int_0^{\frac{16R^2}{\mu^2 n}\Delta^2} \mathbb{P}\left[\|\bar{\theta}_n - \theta_*\|^2 \ge u\right] du$$

$$+ \int_{\frac{16R^2}{\mu^2 n}\Delta^2}^{\frac{16R^2}{\mu^2 n}\left(\frac{\mu\sqrt{n}}{4R^2} + \Delta\right)^2} \mathbb{P}\left[\|\bar{\theta}_n - \theta_*\|^2 \ge u\right] du$$

$$+ \int_{\frac{16R^2}{\mu^2 n}\left(\frac{\mu\sqrt{n}}{4R^2} + \Delta\right)^2}^{\infty} \mathbb{P}\left[\|\bar{\theta}_n - \theta_*\|^2 \ge u\right] du.$$
(22)

The first term in Eq. (22) is simply bounded by bounding the tail bound by one (like in the previous section): $\int_{0}^{\frac{16R^2}{\mu^2 n}\Delta^2} \mathbb{P}\left[\|\bar{\theta}_n - \theta_*\|^2 \ge u\right] du \leqslant \frac{16R^2}{\mu^2 n}\Delta^2$. The last integral in Eq. (22) may be bounded as follows:

$$\begin{split} & \int_{\frac{16R^2}{\mu^2 n}}^{\infty} \left(\frac{\mu\sqrt{n}}{4R^2} + \Delta\right)^2 \mathbb{P}\left[\|\bar{\theta}_n - \theta_*\|^2 \ge u\right] du \\ &= \mathbb{E}\left[1_{\|\bar{\theta}_n - \theta_*\|^2 \ge \frac{16R^2}{\mu^2 n} \left(\frac{\mu\sqrt{n}}{4R^2} + \Delta\right)^2} \|\bar{\theta}_n - \theta_*\|^2\right] \\ &\leqslant \mathbb{P}\left[\|\bar{\theta}_n - \theta_*\|^2 \ge \frac{16R^2}{\mu^2 n} \left(\frac{\mu\sqrt{n}}{4R^2} + \Delta\right)^2\right]^{1/2} \left[\mathbb{E}\left(\|\bar{\theta}_n - \theta_*\|^4\right)\right]^{1/2} \\ &\text{ using Cauchy-Schwarz inequality,} \\ &\leqslant \mathbb{P}\left[\|\bar{\theta}_n - \theta_*\|^2 \ge \frac{16R^2}{\mu^2 n} \left(\frac{\mu\sqrt{n}}{4R^2} + \Delta\right)^2\right]^{1/2} \left(\|\theta_0 - \theta_*\|^2 + 9\gamma^2 nR^2\right) \text{ using Prop. 3.} \end{split}$$

Moreover, if we denote by t_0 the largest solution of $10\sqrt{t_0} + 20\Box t_0 = \frac{\mu\sqrt{n}}{4R^2}$, we have:

$$\sqrt{t_0} = \frac{-10 + \sqrt{100 + 20\Box \frac{\mu\sqrt{n}}{R}}}{40\Box} = \frac{-10 + 10\sqrt{1 + 20\Box \frac{\mu\sqrt{n}}{100R}}}{40\Box}$$

$$\ge \frac{9}{40\Box}\sqrt{20\Box \frac{\mu\sqrt{n}}{100R}},$$

as soon as $20\Box \frac{\mu\sqrt{n}}{100R} \ge 100$, since if $q \ge 100$, $-1 + \sqrt{1+q} \le \frac{9}{10}\sqrt{q}$. This implies that

$$\begin{split} & \int_{\frac{16R^2}{\mu^2 n} \left(\frac{\mu\sqrt{n}}{4R^2} + \Delta\right)^2}^{\infty} \mathbb{P}\left[\|\bar{\theta}_n - \theta_*\|^2 \ge u \right] du \\ \leqslant & \left[4 \exp(-t_0) \right]^{1/2} \left(\|\theta_0 - \theta_*\|^2 + 9\gamma^2 nR^2 \right) \\ \leqslant & \frac{9}{2t_0^2} \left(\|\theta_0 - \theta_*\|^2 + 9\gamma^2 nR^2 \right) \text{ using } \exp(-\alpha) \leqslant \frac{9}{16\alpha^2} \text{ for all } \alpha > 0, \\ \leqslant & \frac{9}{2} \frac{40^4 \Box^4 100^2 R^4}{9^4 20^2 \Box^2 \mu^2 n} \left[\frac{9}{4} \Box^2 / R^2 + \frac{\gamma\sqrt{n}}{3} \Delta \right] \\ \leqslant & 686 \times 64 \frac{\Box^2 R^2}{\mu^2 n} \left[\frac{9}{4} \Box^2 + \frac{1}{6} \Box \Delta \right]. \end{split}$$

The second term in Eq. (22) is bounded exactly like in Appendix F.2, leading to:

$$\begin{split} &\int_{\Delta^2 \frac{16R^2}{\mu^2 n}}^{\frac{16R^2}{\mu^2 n}} \mathbb{P}\big[\|\bar{\theta}_n - \theta_*\|^2 \geqslant u\big] du \\ \leqslant &\int_{0}^{\infty} 4e^{-t} d\bigg(\frac{16R^2}{\mu^2 n} \bigg[10\sqrt{t} + 20\Box t + \Delta\bigg]^2\bigg) \\ \leqslant &\frac{64R^2}{\mu^2 n} \int_{0}^{\infty} e^{-t} \bigg(100 + 400\Box^2 2t + 400\Box\frac{3}{2}t^{1/2} + 20\triangle\frac{1}{2}t^{-1/2} + 40\triangle\Box\bigg) dt \\ \leqslant &\frac{64R^2}{\mu^2 n} \bigg(100\Gamma(1) + 400\Box^2 2\Gamma(2) + 400\Box\frac{3}{2}\Gamma(3/2) + 20\triangle\frac{1}{2}\Gamma(1/2) + 40\triangle\Box\Gamma(1)\bigg) dt \\ \leqslant &\frac{64R^2}{\mu^2 n} \bigg(100 + 400\Box^2 2 + 400\Box\frac{3}{2}\frac{1}{2}\sqrt{\pi} + 20\triangle\frac{1}{2}\sqrt{\pi} + 40\triangle\Box\bigg). \end{split}$$

We can now put all elements together to obtain, from Eq. (22):

$$\begin{split} & \mathbb{E}\|\bar{\theta}_{n} - \theta_{*}\|^{2} \\ \leqslant \quad \frac{64R^{2}}{\mu^{2}n} \bigg(100 + 400\Box^{2}2 + 400\Box\frac{3}{2}\frac{1}{2}\sqrt{\pi} + 20\triangle\frac{1}{2}\sqrt{\pi} + 40\triangle\Box\bigg) \\ & \quad + \frac{16R^{2}}{\mu^{2}n}\triangle^{2} + 686 \times 64\frac{\Box^{2}R^{2}}{\mu^{2}n}\bigg[\frac{9}{4}\Box^{2} + \frac{1}{6}\Box\triangle\bigg] \\ \leqslant \quad \frac{64R^{2}}{n\mu^{2}}\bigg[\frac{1}{4}\triangle^{2} + 100 + 800\Box^{2} + 532\Box + 32\triangle + 40\triangle\Box + 686\frac{9}{4}\Box^{4} + 686\frac{\triangle\Box^{3}}{6}\bigg]. \end{split}$$

BACH

For
$$\gamma = \frac{1}{2R^2\sqrt{N}}$$
, with $\alpha = R \|\theta_0 - \theta_*\|$, $\Box = 1$ and $\triangle = 6\alpha^2 + 6\alpha$, we get

$$\mathbb{E}\|\bar{\theta}_N - \theta_*\|^2 \leqslant \frac{8R^2}{N\mu^2} \Big[2\triangle^2 + 8\triangle(32 + 40 + 115) + 8(100 + 800 + 532 + 1544) \Big]$$

$$\leqslant \frac{8R^2}{N\mu^2} \Big[2\triangle^2 + 1496\triangle + 23808 \Big]$$

$$\leqslant \frac{8R^2}{N\mu^2} \Big[72\alpha^4 + 144\alpha^3 + 72\alpha^2 + 1496 \times 6\alpha^2 + 1496 \times 6\alpha + 23808 \Big]$$

$$\leqslant \frac{R^2}{N\mu^2} \Big[1296\alpha^4 + 18144\alpha^3 + 95256\alpha^2 + 222264\alpha + 194481 \Big] = \frac{R^2}{N\mu^2} (6\alpha + 21)^4$$

The previous bound is valid as long as $\frac{\mu\sqrt{N}}{R} \ge \frac{10000}{20} = 500$. If it is not satisfied, then Prop. 1 shows that it is still valid.

Acknowledgements

This work was supported by the European Research Council (SIERRA Project).

References

- A. Agarwal, P. L. Bartlett, P. Ravikumar, and M. J. Wainwright. Information-theoretic lower bounds on the oracle complexity of stochastic convex optimization. *IEEE Transactions on Information Theory*, 58(5):3235–3249, 2012.
- F. Bach. Self-concordant analysis for logistic regression. *Electronic Journal of Statistics*, 4:384–414, 2010.
- F. Bach and E. Moulines. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In Advances in Neural Information Processing Systems (NIPS), 2011.
- F. Bach and E. Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate O(1/n). Technical Report 00831977, HAL, 2013. To appear in Advances in Neural Information Processing Systems (NIPS).
- A. Bordes, S. Ertekin, J. Weston, and L. Bottou. Fast kernel classifiers with online and active learning. *Journal of Machine Learning Research*, 6:1579–1619, 2005.
- L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In Advances in Neural Information Processing Systems (NIPS), 2008.
- L. Bottou and Y. Le Cun. On-line learning for very large data sets. Applied Stochastic Models in Business and Industry, 21(2):137–151, 2005.
- S. Boucheron, G. Lugosi, and P. Massart. Concentration Inequalities: A Nonasymptotic Theory of Independence. Oxford University Press, 2013.

- M. N. Broadie, D. M. Cicek, and A. Zeevi. General bounds and finite-time improvement for stochastic approximation algorithms. Technical report, Columbia University, 2009.
- J. Duchi and Y. Singer. Efficient online and batch learning using forward backward splitting. *Journal of Machine Learning Research*, 10:2899–2934, 2009.
- V. Fabian. On asymptotic normality in stochastic approximation. *The Annals of Mathematical Statistics*, 39(4):1327–1332, 1968.
- E. Hazan and S. Kale. Beyond the regret minimization barrier: an optimal algorithm for stochastic strongly-convex optimization. In *Proceedings of the Conference on Learning Theory (COLT)*, 2001.
- A. Juditsky and Y. Nesterov. Primal-dual subgradient methods for minimizing uniformly convex functions. Technical Report 00508933, HAL, 2010.
- S. M. Kakade and A. Tewari. On the generalization ability of online strongly convex programming algorithms. In Advances in Neural Information Processing Systems (NIPS), 2009.
- H. J. Kushner and G. G. Yin. *Stochastic approximation and recursive algorithms and applications*. Springer-Verlag, second edition, 2003.
- S. Lacoste-Julien, M. Schmidt, and F. Bach. A simpler approach to obtaining an O(1/t) convergence rate for projected stochastic subgradient descent. Technical Report 1212.2002, ArXiv, 2012.
- J. Lafferty, A. McCallum, and F. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2001.
- G. Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*, 133(1-2):365–397, 2012.
- N. Le Roux, M. Schmidt, and F. Bach. A stochastic gradient method with an exponential convergence rate for strongly-convex optimization with finite training sets. In *Advances in Neural Information Processing Systems (NIPS)*, 2012.
- H. B. McMahan and M. Streeter. Open problem: Better bounds for online logistic regression. In *COLT/ICML Joint Open Problem Session*, 2012.
- A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. SIAM Journal on Optimization, 19(4):1574–1609, 2009.
- A. S. Nemirovsky and D. B. Yudin. Problem complexity and method efficiency in optimization. Wiley & Sons, 1983.
- Y. Nesterov. *Introductory lectures on convex optimization: a basic course*. Kluwer Academic Publishers, 2004.
- Y. Nesterov and A. Nemirovskii. *Interior-point polynomial algorithms in convex programming*. SIAM studies in Applied Mathematics, 1994.

- Y. Nesterov and J. P. Vial. Confidence level solutions for stochastic programming. *Automatica*, 44 (6):1559–1568, 2008.
- I. Pinelis. Optimum bounds for the distributions of martingales in banach spaces. *The Annals of Probability*, 22(4):1679–1706, 1994.
- B. T. Polyak and A. B. Juditsky. Acceleration of stochastic approximation by averaging. SIAM Journal on Control and Optimization, 30(4):838–855, 1992.
- D. Ruppert. Efficient estimations from a slowly convergent Robbins-Monro process. Technical Report 781, Cornell University Operations Research and Industrial Engineering, 1988.
- B. Schölkopf and A. J. Smola. Learning with Kernels. MIT Press, 2001.
- S. Shalev-Shwartz and N. Srebro. SVM optimization: inverse dependence on training set size. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2008.
- S. Shalev-Shwartz, Y. Singer, and N. Srebro. Pegasos: Primal estimated sub-gradient solver for svm. In Proceedings of the International Conference on Machine Learning (ICML), 2007.
- S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Stochastic convex optimization. In Proceedings of the Conference on Learning Theory (COLT), 2009.
- J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.
- A. W. Van der Vaart. Asymptotic Statistics. Cambridge Univ. Press, 1998.
- Z. Wang, K. Crammer, and S. Vucetic. Breaking the curse of kernelization: Budgeted stochastic gradient descent for large-scale SVM training. *Journal of Machine Learning Research*, 13:3103– 3131, 2012.
- L. Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *Journal of Machine Learning Research*, 9:2543–2596, 2010.