
Harder, Better, Faster, Stronger Convergence Rates for Least-Squares Regression

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We consider the optimization of a quadratic objective function whose gradients
2 are only accessible through a stochastic oracle. We present the first algorithm
3 that achieves jointly the optimal prediction error rates for least-squares regression,
4 both in terms of forgetting of initial conditions and in terms of dependence on the
5 noise and dimension of the problem, and prove dimensionless and tighter rates for
6 a regularized version of this algorithm.

7 1 Introduction

8 Many supervised machine learning problems are naturally cast as the minimization of a smooth func-
9 tion defined on a Euclidean space. This includes least-squares regression, logistic regression (see,
10 e.g., Hastie et al., 2009) or generalized linear models (McCullagh and Nelder, 1989). While small
11 problems with few or low-dimensional input features may be solved precisely by many potential
12 optimization algorithms (e.g., Newton method), large-scale problems with many high-dimensional
13 features are typically solved with simple gradient-based iterative techniques whose per-iteration cost
14 is small.

15 In this paper, we consider a quadratic objective function f whose gradients are only accessible
16 through a stochastic oracle that returns the gradient at any given point plus a zero-mean finite vari-
17 ance random error. In this stochastic approximation framework (Robbins and Monro, 1951), it is
18 known that two quantities dictate the behavior of various algorithms, namely the covariance ma-
19 trix V of the noise in the gradients, and the deviation $\theta_0 - \theta_*$ between the initial point of the
20 algorithm θ_0 and any of the global minimizer θ_* of f . This leads to a “bias/variance” decomposi-
21 tion (Bach and Moulines, 2013; Hsu et al., 2014) of the performance of most algorithms as the sum
22 of two terms: (a) the bias term characterizes how fast initial conditions are forgotten and thus is
23 increasing in a well-chosen norm of $\theta_0 - \theta_*$; while (b) the variance term characterizes the effect of
24 the noise in the gradients, independently of the starting point, and with a term that is increasing in
25 the covariance of the noise.

26 For quadratic functions with (a) a noise covariance matrix V which is proportional (with constant
27 σ^2) to the Hessian of f (a situation which corresponds to least-squares regression) and (b) an ini-
28 tial point characterized by the norm $\|\theta_0 - \theta_*\|^2$, the optimal bias and variance terms are known
29 *separately*. On the one hand, the optimal bias term after n iterations is proportional to $\frac{L\|\theta_0 - \theta_*\|^2}{n^2}$,
30 where L is the largest eigenvalue of the Hessian of f . This rate is achieved by accelerated gradient
31 descent (Nesterov, 1983, 2004), and is known to be optimal if the number of iterations n is less than
32 the dimension d of the underlying predictors, but the algorithm is not robust to random or determin-
33 istic noise in the gradients (d’Aspremont, 2008; Schmidt et al., 2011; Devolder et al., 2014). On the
34 other hand, the optimal variance term is proportional to $\frac{\sigma^2 d}{n}$ (Tsybakov, 2003); it is known to be
35 achieved by averaged gradient descent (Bach and Moulines, 2013), for which the bias term only
36 achieves $\frac{L\|\theta_0 - \theta_*\|^2}{n}$ instead of $\frac{L\|\theta_0 - \theta_*\|^2}{n^2}$.

37 Our first contribution in this paper is to analyze in Section 3 averaged *accelerated* gradient descent,
 38 showing that it attains optimal rates for *both the variance and the bias terms*. It shows that averaging
 39 is beneficial for accelerated techniques and provides a provable robustness to noise.

40 While optimal when measuring performance in terms of the dimension d and the initial distance to
 41 optimum $\|\theta_0 - \theta_*\|^2$, these rates are not adapted in many situations where either d is larger than the
 42 number of iterations n (i.e., the number of observations for regular stochastic gradient descent) or
 43 $L\|\theta_0 - \theta_*\|^2$ is much larger than n^2 . Our second contribution is to provide in Section 4 an analysis
 44 of a new algorithm (based on some additional regularization) that can adapt our bounds to finer
 45 assumptions on $\theta_0 - \theta_*$ and the Hessian of the problem, leading in particular to dimension-free
 46 quantities that can thus be extended to the Hilbert space setting (in particular for non-parametric
 47 estimation).

48 2 Least-Squares Regression

49 **Statistical Assumptions.** We consider the following general setting: \mathcal{H} is a d -dimensional Eu-
 50 clidean space with $d \geq 1$, the observations $(x_n, y_n) \in \mathcal{H} \times \mathbb{R}$, $n \geq 1$, are independent and iden-
 51 tically distributed (i.i.d.), and such that $\mathbb{E}\|x_n\|^2$ and $\mathbb{E}y_n^2$ are finite. We consider the *least-squares*
 52 *regression* problem which is the minimization of the quadratic function $f(\theta) = \frac{1}{2}\mathbb{E}(\langle x_n, \theta \rangle - y_n)^2$.

53 *Covariance matrix:* We denote by $\Sigma = \mathbb{E}(x_n \otimes x_n) \in \mathbb{R}^{d \times d}$ the population covariance matrix,
 54 which is the Hessian of f at all points. Without loss of generality, we can assume Σ invertible.
 55 This implies that all eigenvalues of Σ are strictly positive (but they may be arbitrarily small). We
 56 assume there exists $R > 0$ such that $\mathbb{E}\|x_n\|^2 x_n \otimes x_n \preceq R^2 \Sigma$ where $A \preceq B$ means that $B - A$
 57 is positive semi-definite. This assumption is satisfied, for example, for least-square regression with
 58 almost surely bounded data.

59 *Eigenvalue decay:* Most convergence bounds depend on the dimension d of \mathcal{H} . However it is pos-
 60 sible to derive dimension-free and often tighter convergence rates by considering bounds depending
 61 on the value $\text{tr} \Sigma^b$ for $b \in [0, 1]$. Given b , if we consider the eigenvalues of Σ ordered in decreasing
 62 order, which we denote by s_i , then $\text{tr} \Sigma^b = \sum_i s_i^b$, and the eigenvalues decay. For b going to 0 then
 63 $\text{tr} \Sigma^b$ tends to d and we are back in the classical low-dimensional case. When $b = 1$, we simply get
 64 $\text{tr} \Sigma = \mathbb{E}\|x_n\|^2$, which will correspond to the weakest assumption in our context.

65 *Optimal predictor:* The regression function $f(\theta) = \frac{1}{2}\mathbb{E}(\langle x_n, \theta \rangle - y_n)^2$ always admits a global
 66 minimum $\theta_* = \Sigma^{-1}\mathbb{E}(y_n x_n)$. When initializing algorithms at $\theta_0 = 0$ or regularizing by the squared
 67 norm, rates of convergence generally depend on $\|\theta_*\|$, a quantity which could be arbitrarily large.
 68 However there exists a systematic upper-bound $\|\Sigma^{\frac{1}{2}}\theta_*\| \leq 2\sqrt{\mathbb{E}y_n^2}$. This leads naturally to the
 69 consideration of convergence bounds depending on $\|\Sigma^{r/2}\theta_*\|$ for $r \leq 1$.

70 *Noise:* We denote by $\varepsilon_n = y_n - \langle \theta_*, x_n \rangle$ the residual for which we have $\mathbb{E}[\varepsilon_n x_n] = 0$. Although
 71 we do not have $\mathbb{E}[\varepsilon_n | x_n] = 0$ in general unless the model is well-specified, we assume the noise to
 72 be a structured process such that there exists $\sigma > 0$ with $\mathbb{E}[\varepsilon_n^2 x_n \otimes x_n] \preceq \sigma^2 \Sigma$. This assumption is
 73 satisfied for example for data almost surely bounded or when the model is well-specified.

74 **Averaged Gradient Methods and Acceleration.** We focus in this paper on stochastic gradient
 75 methods with acceleration for a quadratic function regularized by $\frac{\lambda}{2}\|\theta - \theta_0\|^2$. The regularization
 76 will be useful when deriving tighter convergence rates in Section 4, and it has the additional benefit
 77 of making the problem λ -strongly-convex.

78 Accelerated stochastic gradient descent is defined by an iterative system with two parameters
 79 (θ_n, ν_n) starting from $\theta_0 = \nu_0 \in \mathcal{H}$, and satisfying for $n \geq 1$,

$$\begin{aligned} \theta_n &= \nu_{n-1} - \gamma f'_n(\nu_{n-1}) - \gamma \lambda (\nu_{n-1} - \theta_0) \\ \nu_n &= \theta_n + \delta (\theta_n - \theta_{n-1}), \end{aligned} \tag{1}$$

80 with $\gamma, \delta \in \mathbb{R}^2$ and $f'_n(\theta_{n-1})$ an unbiased estimate on the gradient $f(\theta)$.

81 The *momentum* coefficient $\delta \in \mathbb{R}$ is chosen to accelerate the convergence rate (Nesterov, 1983;
 82 Beck and Teboulle, 2009) and has its roots in the heavy-ball algorithm from Polyak (1964). We
 83 especially concentrate here, following Polyak and Juditsky (1992), on the average of the sequence
 84 $\bar{\theta}_n = \frac{1}{n+1} \sum_{i=0}^n \theta_i$,

85 **Stochastic Oracles on the Gradient.** Let $(\mathcal{F}_n)_{n \geq 0}$ be the increasing family of σ -fields that are
 86 generated by all variables (x_i, y_i) for $i \leq n$. The oracle we consider is the sum of the true gradient
 87 $f'(\theta)$ and an independent zero-mean noise that does not depend on θ^1 . Consequently it is of the
 88 form $f'_n(\theta) = f'(\theta) - \xi_n$ where the noise process ξ_n is \mathcal{F}_n -measurable with $\mathbb{E}[\xi_n | \mathcal{F}_{n-1}] = 0$ and
 89 $\mathbb{E}[\|\xi_n\|^2]$ is finite. Furthermore we also assume that there exists $\tau \in \mathbb{R}$ such that $\mathbb{E}[\xi_n \otimes \xi_n] \preceq \tau^2 \Sigma$,
 90 that is, the noise has a particular structure adapted to least-squares regression.

91 3 Accelerated Stochastic Averaged Gradient Descent

92 We study the convergence of averaged *accelerated* stochastic gradient descent defined by Eq. (1) for
 93 $\lambda = 0$ and $\delta = 1$. It can be rewritten for the quadratic function f as a second-order iterative system
 94 with constant coefficients: $\theta_n = [I - \gamma \Sigma](2\theta_{n-1} - \theta_{n-2}) + \gamma y_n x_n$.

95 **Theorem 1** For any constant step-size γ , such that $\gamma \Sigma \preceq I$,

$$\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leq 36 \frac{\|\theta_0 - \theta_*\|^2}{\gamma(n+1)^2} + 8 \frac{\tau^2 d}{n+1}. \quad (2)$$

96 We can make the following observations:

- 97 • The first bound $\frac{1}{\gamma n^2} \|\theta_0 - \theta_*\|^2$ in Eq. (2) corresponds to the usual accelerated rate. It has
 98 been shown by Nesterov (2004) to be the optimal rate of convergence for optimizing a
 99 quadratic function with a first-order method that can access only to sequences of gradients
 100 when $n \leq d$. We recover by averaging an algorithm dedicated to strongly-convex function
 101 the traditional convergence rate for non-strongly convex functions.
- 102 • The second bound in Eq. (2) matches the optimal statistical performance $\frac{\tau^2 d}{n}$ over all
 103 estimators in \mathcal{H} (Tsybakov, 2008) even without computational limits, in the sense that
 104 no estimator that uses the same information can improve upon this rate. Accordingly
 105 this algorithm achieves joint bias/variance optimality (when measured in terms of τ^2 and
 106 $\|\theta_0 - \theta_*\|^2$).
- 107 • We have the same rate of convergence for the bias when compared to the regular Nesterov
 108 acceleration without averaging studied by Flammarion and Bach (2015), which cor-
 109 responds to choosing $\delta_n = 1 - 2/n$ for all n . However if the problem is μ -strongly convex,
 110 this latter was shown to also converge at the linear rate $O((1 - \gamma\mu)^n)$ and thus is adaptive
 111 to hidden strong-convexity (since the algorithm does not need to know μ to run), thus ends
 112 up converging faster than the rate $1/n^2$. This is confirmed in our experiments in Section 5.
- 113 • Overall, the bias term is improved whereas the variance term is not degraded and accelera-
 114 tion is thus robust to noise in the gradients. Thereby, while second-order iterative methods
 115 for optimizing quadratic functions in the singular case, such as conjugate gradient (Polyak,
 116 1987, Section 6.1) are notoriously highly sensitive to noise, we are able to propose a version
 117 which is robust to stochastic noise.

118 4 Tighter Convergence Rates

119 We have seen in Corollary 1 above that the averaged accelerated gradient algorithm matches the
 120 lower bounds $\tau^2 d/n$ and $\frac{L}{n^2} \|\theta_0 - \theta_*\|^2$ for the prediction error. However the algorithm performs
 121 better in almost all cases except the worst-case scenarios corresponding to the lower bounds. For
 122 example the algorithm may still predict well when the dimension d is much bigger than n . Similarly
 123 the norm of the optimal predictor $\|\theta_*\|^2$ may be huge and the prediction still good, as gradient
 124 algorithms happen to be adaptive to the difficulty of the problem: indeed, if the problem is simpler,
 125 the convergence rate of the gradient algorithm will be improved. In this section, we provide such a
 126 theoretical guarantee.

127 We study the convergence of averaged *accelerated* stochastic gradient descent defined by Eq. (1) for
 128 $\lambda = (\gamma(n+1)^2)^{-1}$ and $\delta \in [1 - \frac{2}{n+2}, 1]$. We have the following theorem:

¹this is different from the oracle usually considered in stochastic approximation (see Bach and Moulines (2013); Dieuleveut and Bach (2015)).

129 **Theorem 2** For any constant step-size γ , such that $\gamma(\Sigma + \lambda I) \preceq I$,

$$\mathbb{E}f(\bar{\theta}_n) - f(\theta_*) \leq \min_{r \in [0,1], b \in [0,1]} \left[74 \frac{\|\Sigma^{r/2}(\theta_0 - \theta_*)\|^2}{\gamma^{1-r}(n+1)^{2(1-r)}} + 8 \frac{\tau^2 \gamma^b \text{tr}(\Sigma^b)}{(n+1)^{1-2b}} \right].$$

130 We can make the following observations:

- 131 • The algorithm is independent of r and b , thus all the bounds for different values of (r, b)
- 132 are valid. This is a strong property of the algorithm, which is indeed adaptative to the
- 133 regularity and the effective dimension of the problem (once γ is chosen). In situations in
- 134 which either d is larger than n or $L\|\theta_0 - \theta_*\|^2$ is larger than n^2 , the algorithm can still enjoy
- 135 good convergence properties, by adapting to the best values of b and r .
- 136 • For $b = 0$ we recover the variance term of Corollary 1, but for $b > 0$ and fast decays of
- 137 eigenvalues of Σ , the bound may be much smaller; note that we lose in the dependency in
- 138 n , but typically, for large d , this can be advantageous.
- 139 • With r, b well chosen, we recover the optimal rate for non-parametric regression
- 140 (Caponnetto and De Vito, 2007).

141 5 Experiments

142 We illustrate now our theoretical results on synthetic examples. For $d = 25$ we consider normally

143 distributed inputs x_n with random covariance matrix Σ which has eigenvalues $1/i^3$, for $i = 1, \dots, d$,

144 and random optimum θ_* and starting point θ_0 such that $\|\theta_0 - \theta_*\| = 1$. The outputs y_n are generated

145 from a linear function with homoscedastic noise with unit signal to noise-ratio ($\sigma^2 = 1$), we take

146 $R^2 = \text{tr} \Sigma$ the average radius of the data and a step-size $\gamma = 1/R^2$ and $\lambda = 0$. The additive noise

147 oracle is used. We show results averaged over 10 replications.

148 We compare the performance of averaged SGD (AvSGD), usual Nesterov acceleration for convex

149 functions (AccSGD) and our novel averaged accelerated SGD (AvAccSGD)², on two different prob-

150 lems: one deterministic ($\|\theta_0 - \theta_*\| = 1, \sigma^2 = 0$) which will illustrate how the bias term behaves,

151 and one purely stochastic ($\|\theta_0 - \theta_*\| = 0, \sigma^2 = 1$) which will illustrate how the variance term

behaves. For the bias (left plot of Figure 1), AvSGD converges at speed $O(1/n)$, while AvAccSGD

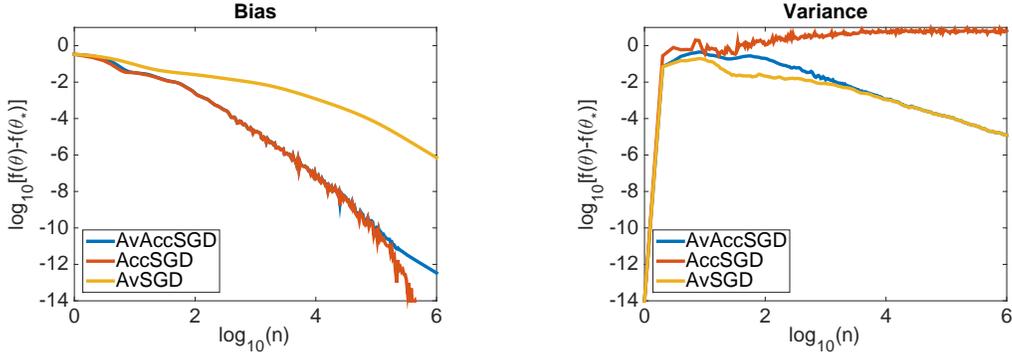


Figure 1: Synthetic problem ($d = 25$) and $\gamma = 1/R^2$. Left: Bias. Right: Variance.

152 and AccSGD converge both at speed $O(1/n^2)$. However, as mentioned in the observations follow-

153 ing Theorem 1, AccSGD takes advantage of the hidden strong convexity of the quadratic function

154 and starts converging linearly at the end. For the variance (right plot of Figure 1), AccSGD is not

155 converging to the optimum and keeps oscillating whereas AvSGD and AvAccSGD both converge

156 to the optimum at a speed $O(1/n)$. However AvSGD remains slightly faster in the beginning.

158 Note that for small n , or when the bias $L\|\theta_0 - \theta_*\|^2/n^2$ is much bigger than the variance $\sigma^2 d/n$, the

159 bias may have a stronger effect, although asymptotically, the variance always dominates. It is thus

160 essential to have an algorithm which is optimal in both regimes; this is achieved by AvAccSGD.

²which is not the averaging of AccSGD because the momentum term is proportional to $1 - 3/n$ for AccSGD instead of being equal to 1 for AvAccSGD.

161 **References**

- 162 Bach, F. and Moulines, E. (2013). Non-strongly-convex smooth stochastic approximation with
163 convergence rate $O(1/n)$. In *Advances in Neural Information Processing Systems*.
- 164 Beck, A. and Teboulle, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear
165 inverse problems. *SIAM J. Imaging Sci.*, 2(1):183–202.
- 166 Caponnetto, A. and De Vito, E. (2007). Optimal rates for the regularized least-squares algorithm.
167 *Foundations of Computational Mathematics*, 7(3):331–368.
- 168 d’Aspremont, A. (2008). Smooth optimization with approximate gradient. *SIAM J. Optim.*,
169 19(3):1171–1183.
- 170 Devolder, O., Glineur, F., and Nesterov, Y. (2014). First-order methods of smooth convex optimiza-
171 tion with inexact oracle. *Math. Program.*, 146(1-2, Ser. A):37–75.
- 172 Dieuleveut, A. and Bach, F. (2015). Non-parametric stochastic approximation with large step sizes.
173 *Annals of Statistics*. To appear.
- 174 Flammarion, N. and Bach, F. (2015). From averaging to acceleration, there is only a step-size. In
175 *Proceedings of the International Conference on Learning Theory (COLT)*.
- 176 Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The Elements of Statistical Learning*. Springer
177 Series in Statistics. Springer, second edition.
- 178 Hsu, D., Kakade, S. M., and Zhang, T. (2014). Random design analysis of ridge regression. *Foun-
179 dations of Computational Mathematics*, 14(3):569–600.
- 180 McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*. Monographs on Statistics and
181 Applied Probability. Chapman & Hall, London, second edition.
- 182 Nesterov, Y. (1983). A method of solving a convex programming problem with convergence rate
183 $O(1/k^2)$. *Soviet Mathematics Doklady*, 27(2):372–376.
- 184 Nesterov, Y. (2004). *Introductory Lectures on Convex Optimization*, volume 87 of *Applied Optimiza-
185 tion*. Kluwer Academic Publishers, Boston, MA.
- 186 Polyak, B. T. (1964). Some methods of speeding up the convergence of iteration methods. *{USSR}
187 Computational Mathematics and Mathematical Physics*, 4(5):1–17.
- 188 Polyak, B. T. (1987). *Introduction to Optimization*. Translations Series in Mathematics and Engi-
189 neering. Optimization Software, Inc., Publications Division, New York.
- 190 Polyak, B. T. and Juditsky, A. B. (1992). Acceleration of stochastic approximation by averaging.
191 *SIAM J. Control Optim.*, 30(4):838–855.
- 192 Robbins, H. and Monroe, S. (1951). A stochastic approximation method. *The Annals of mathematical
193 Statistics*, 22(3):400–407.
- 194 Schmidt, M., Le Roux, N., and Bach, F. (2011). Convergence rates of inexact proximal-gradient
195 methods for convex optimization. In *Advances in Neural Information Processing Systems*.
- 196 Tsybakov, A. B. (2003). Optimal rates of aggregation. In *Proceedings of the Annual Conference on
197 Computational Learning Theory*.
- 198 Tsybakov, A. B. (2008). *Introduction to Nonparametric Estimation*. Springer.