

An Accelerated Method for Derivative-Free Smooth Stochastic Convex Optimization

Pavel Dvurechensky

Weierstrass Institute for Applied Analysis and Stochastics

PAVEL.DVURECHENSKY@WIAS-BERLIN.DE

Alexander Gasnikov

Moscow Institute of Physics and Technology

GASNIKOV@YANDEX.RU

Eduard Gorbunov

Moscow Institute of Physics and Technology

ED-GORBUNOV@YANDEX.RU

Abstract

We consider an unconstrained problem of minimization of a smooth convex function which is only available through noisy observations of its values, the noise consisting of two parts. Similar to stochastic optimization problems, the first part is of a stochastic nature. On the opposite, the second part is an additive noise of an unknown nature, but bounded in the absolute value. In the two-point feedback setting, i.e. when pairs of function values are available, we propose an accelerated derivative-free algorithm together with its complexity analysis. The complexity bound of our derivative-free algorithm is only by a factor of \sqrt{n} larger than the bound for accelerated gradient-based algorithms, where n is the dimension of the decision variable. We also propose a non-accelerated derivative-free algorithm with a complexity bound similar to the stochastic-gradient-based algorithm, that is, our bound does not have any dimension-dependent factor. Interestingly, if the solution of the problem is sparse, for both our algorithms, we obtain better complexity bound if the algorithm uses a 1-norm proximal setup, rather than the Euclidean proximal setup, which is a standard choice for unconstrained problems.

Keywords: Derivative-Free Optimization, Stochastic Convex Optimization, Smoothness, Acceleration

1. Introduction

Derivative-free optimization [Rosenbrock \(1960\)](#); [Brent \(1973\)](#); [Spall \(2003\)](#) is one of the oldest areas in optimization, which constantly attracts attention of the learning community, mostly in connection to the online learning in the bandit setup [Bubeck and Cesa-Bianchi \(2012\)](#). We study stochastic derivative-free optimization problems in a two-point feedback situation, considered by [Agarwal et al. \(2010\)](#); [Duchi et al. \(2015\)](#); [Shamir \(2017\)](#) in the learning community and by [Nesterov and Spokoiny \(2017\)](#); [Stich et al. \(2011\)](#); [Ghadimi and Lan \(2013\)](#); [Ghadimi et al. \(2016\)](#); [Gasnikov et al. \(2016b\)](#) in the optimization community. Two-point setup allows to prove complexity bounds, which typically coincide with the complexity bounds for gradient-based algorithms up to a small-degree polynomial of n , where n is the dimension of the decision variable. On the contrary, problems with one-point feedback are harder and complexity bounds for such problems either have worse dependence on n , or worse dependence on the desired accuracy of the solution, see [Flaxman et al. \(2005\)](#); [Agarwal et al. \(2011\)](#); [Jamieson et al. \(2012\)](#); [Shamir \(2013\)](#); [Liang et al. \(2014\)](#); [Bach and Perchet \(2016\)](#); [Bubeck et al. \(2017\)](#) and the references therein.

More precisely, we consider the following optimization problem

$$\min_{x \in \mathbb{R}^n} \left\{ f(x) := \mathbb{E}_\xi[F(x, \xi)] = \int_{\mathcal{X}} F(x, \xi) dP(\xi) \right\}, \quad (1)$$

where ξ is a random vector with probability distribution $P(\xi)$, $\xi \in \mathcal{X}$, and for P -almost every $\xi \in \mathcal{X}$, the function $F(x, \xi)$ is closed and convex. Moreover, we assume that, for P almost every ξ , the function $F(x, \xi)$ has gradient $g(x, \xi)$, which is $L(\xi)$ -Lipschitz continuous with respect to the Euclidean norm and $L_2 := \sqrt{\mathbb{E}_\xi L(\xi)^2} < +\infty$. Under this assumptions, $\mathbb{E}_\xi g(x, \xi) = \nabla f(x)$ and f has L_2 -Lipschitz continuous gradient with respect to the Euclidean norm. Also we assume that

$$\mathbb{E}_\xi[\|g(x, \xi) - \nabla f(x)\|_2^2] \leq \sigma^2, \quad (2)$$

where $\|\cdot\|_2$ is the Euclidean norm.

Finally, we assume that an optimization procedure, given a pair of points $(x, y) \in \mathbb{R}^{2n}$, can obtain a pair of noisy stochastic realizations $(f(x, \xi), \tilde{f}(y, \xi))$ of the objective value f , where

$$\tilde{f}(x, \xi) = F(x, \xi) + \eta(x, \xi), \quad |\eta(x, \xi)| \leq \Delta, \quad \forall x \in \mathbb{R}^n, \text{ a.s. in } \xi, \quad (3)$$

and ξ is independently drawn from P .

It is well-known [Lan \(2012\)](#); [Devolder \(2011\)](#); [Dvurechensky and Gasnikov \(2016\)](#) that, if the stochastic approximation $g(x, \xi)$ for the gradient of f is available, an accelerated gradient method has complexity bound $O\left(\max\left\{\sqrt{L_2/\varepsilon}, \sigma^2/\varepsilon^2\right\}\right)$, where ε is the target optimization error. The question, to which we give a positive answer in this paper, is as follows.

Is it possible to solve a stochastic optimization problem with the same ε -dependence in the complexity and only noisy observations of the objective value?

1.1. Related Work

Complexity bounds for derivative-free optimization with exact function observations were obtained in [Nemirovsky and Yudin \(1983\)](#); [Protasov \(1996\)](#). The related work on stochastic problems can be divided in two large groups, namely, a group, considering one-point feedback, and a group, considering two-point feedback. A unified view on these two cases was presented in [Hu et al. \(2016\)](#).

One-point feedback. Strictly speaking, this setup allows to form an approximation for the gradient using two observations of the objective function, but these two observations correspond to two different realizations of the random vector ξ . Most of the authors in this group solve a more general problem of bandit convex optimization and obtain bounds on the so-called regret. It is well known [Cesa-bianchi et al. \(2002\)](#) that a bound on the regret can be converted to a bound on the expected optimization error $\mathbb{E}f(\hat{x}) - f^*$ in stochastic optimization, where f^* is an optimal value of f . To compare the results in the literature, we compare complexity bounds, that is, the number of iterations to achieve the expected optimization error of ε .

In the early work, [Flaxman et al. \(2005\)](#) obtained complexity¹ $\tilde{O}(n^2/\varepsilon^4)$ for convex non-smooth problems and $\tilde{O}(n^2/\varepsilon^3)$ for strongly convex non-smooth problems (see also [Saha and Tewari \(2011\)](#); [Dekel et al. \(2015\)](#); [Gasnikov et al. \(2017\)](#) on some improvements on these sub $1/\varepsilon^2$ rates). [Agarwal et al.](#)

1. \tilde{O} hides polylogarithmic factors $(\ln n)^c$, $c > 0$.

(2011) provide an algorithm with complexity bound $\tilde{O}(n^{32}/\varepsilon^2)$, which was later improved by Liang et al. (2014) to $\tilde{O}(n^{14}/\varepsilon^2)$ for convex functions, and by Belloni et al. (2015) to $\tilde{O}(n^{6.5}/\varepsilon^2)$. Bubeck et al. (2017) conjecture that their algorithm has complexity $\tilde{O}(n^3/\varepsilon^2)$, while the lower bound of Shamir (2013) is $\Omega(n^2/\varepsilon^2)$ even for strongly convex functions.

Smoothness of the objective function allows to obtain better upper bounds. In this case, Jamieson et al. (2012); Hazan and Levy (2014) proved $\tilde{O}(n^3/\varepsilon^2)$ bound for strongly convex problems. Later, Gasnikov et al. (2017) obtained a bound $\tilde{O}(n/\varepsilon^3)$ for convex problems and $\tilde{O}(n^2/\varepsilon^2)$ for strongly convex problems. Bach and Perchet (2016) obtained a bound $\tilde{O}(n^2/\varepsilon^3)$ for convex problems and $\tilde{O}(n^2/\varepsilon^2)$ for strongly convex problems. For the smooth case and both convex and strongly convex problems, Jamieson et al. (2012) proved a lower bound $\Omega(n/\varepsilon^2)$ and Shamir (2013) obtained an $\Omega(n^2/\varepsilon^2)$ lower bound.

Two-point feedback. Non-smooth problem of this type was considered by Nesterov and Spokoiny (2017)², who proved an $O(n^2/\varepsilon^2)$ complexity bound. This bound was improved by Duchi et al. (2015); Gasnikov et al. (2016a, 2017); Shamir (2017); Bayandina et al. (2018) to $\tilde{O}(n^{2/q}R_p^2/\varepsilon^2)$, where $p \in \{1, 2\}$, $\frac{1}{p} + \frac{1}{q} = 1$ and R_p is the radius of the feasible set in the $\|\cdot\|_p$ -norm. For non-smooth μ_p -strongly convex w.r.t. to $\|\cdot\|_p$ -norm problems, Gasnikov et al. (2017); Bayandina et al. (2018) proved a bound $\tilde{O}(n^{2/q}/(\mu_p\varepsilon))$.

Intermediate, partially smooth problems with a restrictive assumption of boundedness of $\mathbb{E}\|g(x, \xi)\|^2$, were considered by Duchi et al. (2015), who proved that a properly modified Mirror Descent algorithm gives a bound $O(n^{2/q}R_p^2/\varepsilon^2)$ for convex problems, improving upon the bound $\tilde{O}(n^2/\varepsilon^2)$ of Agarwal et al. (2010). For convex problems, Agarwal et al. (2010) obtained a bound $\tilde{O}(n^2/\varepsilon)$, which was later extended for μ_p -strongly convex problems and improved to $\tilde{O}(n^{2/q}/(\mu_p\varepsilon))$ in Gasnikov et al. (2017).

In the fully smooth case, without the assumption that $\mathbb{E}\|g(x, \xi)\|^2 < +\infty$, Ghadimi et al. (2016); Ghadimi and Lan (2013) proposed an algorithm with the bound

$$\tilde{O}\left(\max\left\{\frac{nL_2R_2}{\varepsilon}, \frac{n\sigma^2}{\varepsilon^2}\right\}\right)$$

for the Euclidean case.

Deterministic problems. Accelerated derivative-free method for deterministic problems was proposed in Nesterov and Spokoiny (2017) for the Euclidean case with the bound $O(n\sqrt{L_2/\varepsilon})$. A non-accelerated derivative-free method for deterministic problems with additional bounded noise in function values was proposed in Bogolubsky et al. (2016) together with $O(nL_2/\varepsilon)$ bound and application to learning parameter of a parametric PageRank model. Deterministic problems with additional bounded noise in function values were also considered in Dvurechensky et al. (2017), where several accelerated derivative-free methods, including Derivative-Free Block-Coordinate Descent, were proposed and a bound $O(n\sqrt{L/\varepsilon})$ was proved, where L depends on the method and, in some sense, characterizes the average over blocks of coordinates Lipschitz constant of the derivative in the block.

2. We list the references in the order of the date of the first appearance, but not in the order of the date of publication.

1.2. Our Contributions

We solve the problem (1) using two proximal setups [Ben-Tal and Nemirovski \(2015\)](#), characterized by the value³ $p \in \{1, 2\}$ and its conjugate $q \in \{2, \infty\}$, given by the identity $\frac{1}{p} + \frac{1}{q} = 1$. The case $p = 1$ corresponds to the choice of $\|\cdot\|_1$ -norm in \mathbb{R}^n and corresponding prox-function, which is strongly convex with respect to this norm (we provide the details below). The case $p = 2$ corresponds to the choice of the Euclidean $\|\cdot\|_2$ -norm in \mathbb{R}^n and squared Euclidean norm as the prox-function. As our main contribution, we propose an accelerated method for smooth stochastic derivative-free optimization, which we call Accelerated Randomized Derivative-Free Directional Search (ARDFDS). Our method has the complexity bound

$$\tilde{O} \left(\max \left\{ n^{\frac{1}{2} + \frac{1}{q}} \sqrt{\frac{L_2 R_p^2}{\varepsilon}}, \frac{n^{\frac{2}{q}} \sigma^2 R_p^2}{\varepsilon^2} \right\} \right), \quad (4)$$

where R_p characterizes the distance in $\|\cdot\|_p$ -norm between the starting point of the algorithm and a solution to (1). In the Euclidean case $p = q = 2$, the first term in the above bound has better dependence on n, ε, L_2 and R_2 than the bound in [Ghadimi et al. \(2016\)](#); [Ghadimi and Lan \(2013\)](#). Unlike these papers, our bound also covers the non-euclidean case $p = 1, q = \infty$. At the same time, in the second term in (4) we have the same dependence on n as in the case of non-smooth and partially smooth problems in [Duchi et al. \(2015\)](#); [Gasnikov et al. \(2016a, 2017\)](#); [Shamir \(2017\)](#); [Bayandina et al. \(2018\)](#). We also underline that, in the case of (1) having a sparse solution, our bound for $p = 1$ allows to gain a factor of \sqrt{n} in comparison to the Euclidean case $p = 2$. Indeed, if $\|x^*\|_1 \leq O(1) \cdot \|x^*\|_2$ and the starting point is zero, we obtain $\tilde{O} \left(\max \left\{ \sqrt{\frac{n L_2 \|x^*\|_2^2}{\varepsilon}}, \frac{\sigma^2 \|x^*\|_2^2}{\varepsilon^2} \right\} \right)$

for the case $p = 1$ instead of $\tilde{O} \left(\max \left\{ \sqrt{\frac{n^2 L_2 \|x^*\|_2^2}{\varepsilon}}, \frac{n \sigma^2 \|x^*\|_2^2}{\varepsilon^2} \right\} \right)$ for the case $p = 2$. It should also be mentioned that our assumption $\mathbb{E}_\xi L(\xi)^2 < +\infty$ is weaker than the assumption $L(\xi) \leq L_2$ a.s. in ξ , which is used in [Ghadimi et al. \(2016\)](#); [Ghadimi and Lan \(2013\)](#). At the same time, we consider a more general definition of noisy observations (3), while others consider the case $\eta(x, \xi) \equiv 0$.

As our second contribution, we propose a non-accelerated Randomized Derivative-Free Directional Search (RDFDS) method with the complexity bound

$$\tilde{O} \left(\max \left\{ \frac{n^{\frac{2}{q}} L_2 R_p^2}{\varepsilon}, \frac{n^{\frac{2}{q}} \sigma^2 R_p^2}{\varepsilon^2} \right\} \right), \quad (5)$$

where unlike [Ghadimi et al. \(2016\)](#); [Ghadimi and Lan \(2013\)](#) a non-euclidean case $p = 1, q = \infty$ is possible. Interestingly, in this case, we obtain a dimension independent complexity bound despite we use only noisy function value observations.

Note that our results for accelerated and non-accelerated methods are somewhat similar to the finite-sum minimization problems of the form

$$\sum_{i=1}^m f_i(x),$$

3. Strictly speaking, we are able to consider all the intermediate cases $p \in [1, 2]$, but we are not aware of any proximal setup, which is compatible with $p \notin \{1, 2\}$

where f_i are convex smooth functions. For such problems accelerated methods have complexity $\tilde{O}(m + \sqrt{mL/\varepsilon})$ and non-accelerated methods have complexity $\tilde{O}(m + L/\varepsilon)$ (see, e.g. [Allen-Zhu \(2017\)](#) for a nice review on the topic), so acceleration allows to take the square root of the second term but for the price of \sqrt{m} and the two bounds can not be directly compared without additional assumptions on the value of $m\varepsilon$.

2. Algorithms for Stochastic Convex Optimization

2.1. Preliminaries

Proximal setup. Let $p \in [1, 2]$ and $\|x\|_p$ be the p -norm in \mathbb{R}^n defined as

$$\|x\|_p^p = \sum_{i=1}^n |x_i|^p, \quad x \in \mathbb{R}^n$$

and $\|\cdot\|_q$ be its dual, defined by $\|g\|_q = \max_x \{ \langle g, x \rangle, \|x\|_p \leq 1 \}$, where $q \in [2, \infty]$ is the conjugate number to p , given by $\frac{1}{p} + \frac{1}{q} = 1$, and, for $q = \infty$, we define $\|x\|_\infty = \max_{i=1, \dots, n} |x_i|$.

We choose a *prox-function* $d(x)$, which is continuous, convex on \mathbb{R}^n and is 1-strongly convex on \mathbb{R}^n with respect to $\|\cdot\|_p$, i.e., for any $x, y \in \mathbb{R}^n$ $d(y) - d(x) - \langle \nabla d(x), y - x \rangle \geq \frac{1}{2} \|y - x\|_p^2$. Without loss of generality, we assume that $\min_{x \in \mathbb{R}^n} d(x) = 0$. We define also the corresponding *Bregman divergence* $V[z](x) = d(x) - d(z) - \langle \nabla d(z), x - z \rangle$, $x, z \in \mathbb{R}^n$. Note that, by the strong convexity of d ,

$$V[z](x) \geq \frac{1}{2} \|x - z\|_p^2, \quad x, z \in \mathbb{R}^n. \quad (6)$$

For the case $p = 1$, we choose the following prox function [Ben-Tal and Nemirovski \(2015\)](#)

$$d(x) = \frac{en^{(\kappa-1)(2-\kappa)/\kappa} \ln n}{2} \|x\|_\kappa^2, \quad \kappa = 1 + \frac{1}{\ln n}$$

and, for the case $p = 2$, we choose the prox function to be the squared Euclidean norm

$$d(x) = \frac{1}{2} \|x\|_2^2.$$

Main technical lemma. In our proofs of complexity bounds, we rely on the following lemma. The proof is rather technical and is provided in the appendix.

Lemma 1 *Let $e \in RS_2(1)$, i.e be a random vector uniformly distributed on the surface of the unit Euclidean sphere in \mathbb{R}^n , $p \in [1, 2]$ and q be given by $\frac{1}{p} + \frac{1}{q} = 1$. Then, for $n \geq 8$,*

$$\mathbb{E}_e \|e\|_q^2 \leq \rho_n, \quad (7)$$

$$\mathbb{E}_e (\langle s, e \rangle^2 \|e\|_q^2) \leq \frac{6\rho_n}{n} \|s\|_2^2, \quad \forall s \in \mathbb{R}^n. \quad (8)$$

with $\rho_n = \min\{q - 1, 16 \ln n - 8\} n^{\frac{2}{q}-1}$.

Stochastic approximation of the gradient. Based on the noisy observations (3) of the objective value, we form the following stochastic approximation of $\nabla f(x)$

$$\tilde{\nabla}^m f^t(x) = \frac{1}{m} \sum_{i=1}^m \frac{\tilde{f}(x + te, \xi_i) - \tilde{f}(x, \xi_i)}{t} e, \quad (9)$$

where $e \in RS_2(1)$, $\xi_i, i = 1, \dots, m$ are independent realizations of ξ , m is the *batch size*, t is some small positive parameter, which we call *smoothing parameter*.

2.2. Algorithms and Main Theorems

Our Accelerated Randomized Derivative-Free Directional Search (ARDFDS) method is listed as Algorithm 1.

Algorithm 1 Accelerated Randomized Derivative-Free Directional Search (ARDFDS)

Input: x_0 — starting point; N — number of iterations; $m \geq 1$ — batch size; $t > 0$ — smoothing parameter.

Output: point y_N .

- 1: $y_0 \leftarrow x_0, z_0 \leftarrow x_0$.
 - 2: **for** $k = 0, \dots, N - 1$. **do**
 - 3: $\alpha_{k+1} \leftarrow \frac{k+2}{96n^2\rho_n L_2}, \tau_k \leftarrow \frac{1}{48\alpha_{k+1}n^2\rho_n L_2} = \frac{2}{k+2}$.
 - 4: Generate $e_{k+1} \in RS_2(1)$ independently from previous iterations and $\xi_i, i = 1, \dots, m$ — independent realizations of ξ .
 - 5: Calculate $\tilde{\nabla}^m f^t(x_{k+1})$ given in (9).
 - 6: $x_{k+1} \leftarrow \tau_k z_k + (1 - \tau_k) y_k$.
 - 7: $y_{k+1} \leftarrow x_{k+1} - \frac{1}{2L_2} \tilde{\nabla}^m f^t(x_{k+1})$.
 - 8: $z_{k+1} \leftarrow \operatorname{argmin}_{z \in \mathbb{R}^n} \left\{ \alpha_{k+1} n \left\langle \tilde{\nabla}^m f^t(x_{k+1}), z - z_k \right\rangle + V[z_k](z) \right\}$.
 - 9: **end for**
 - 10: **return** y_N
-

Theorem 2 Let ARDFDS be applied to solve problem (1). Then

$$\begin{aligned} \mathbb{E}[f(y_N)] - f(x^*) &\leq \frac{384n^2\rho_n L_2 \Theta_p}{N^2} + \frac{384N}{nL_2} \frac{\sigma^2}{m} \\ &+ \frac{12\sqrt{2n\Theta}}{N^2} \left(\frac{L_2 t}{2} + \frac{2\Delta}{t} \right) + \frac{6N}{L_2} \left(L_2^2 t^2 + \frac{16\Delta^2}{t^2} \right) + \frac{N^2}{24n\rho_n L_2} \left(L_2^2 t^2 + \frac{16\Delta^2}{t^2} \right), \end{aligned} \quad (10)$$

where $\Theta_p = V[z_0](x^*)$ is defined by the chosen proximal setup and the $\mathbb{E}[\cdot] = \mathbb{E}_{e_1, \dots, e_N, \xi_{1,1}, \dots, \xi_{N,m}}[\cdot]$.

The proof is rather technical and is provided in the next subsections. Before we start with the proofs, we give the appropriate choice of the ARDFDS parameters N, m, t and accuracy of the function values evaluation Δ , which we consider to be controlled. These results are obtained first on the choice $t = \sqrt{\Delta/L_2}$, which minimizes $\frac{L_2 t}{2} + \frac{2\Delta}{t}$ and similar expressions in the r.h.s. of (10). Then N, m and Δ are chosen such that the r.h.s. of (10) is smaller than ε . We summarize the obtained values of the algorithm parameters in Table 1 below. The last row represents the total number Nm of function evaluations, which was advertised in (4).

Our Randomized Derivative-Free Directional Search (RDFDS) method is listed as Algorithm 2.

	$p = 1$	$p = 2$
N	$O\left(\sqrt{\frac{n \ln n L_2 \Theta_1}{\varepsilon}}\right)$	$O\left(\sqrt{\frac{n^2 L_2 \Theta_2}{\varepsilon}}\right)$
m	$O\left(\max\left\{1, \frac{\sigma^2 \sqrt{\ln n}}{\varepsilon^{3/2} \sqrt{n}} \cdot \sqrt{\frac{\Theta_1}{L_2}}\right\}\right)$	$O\left(\max\left\{1, \frac{\sigma^2}{\varepsilon^{3/2}} \cdot \sqrt{\frac{\Theta_2}{L_2}}\right\}\right)$
t	$O\left(\min\left\{\frac{\varepsilon^{3/4}}{(n \ln n)^{1/4}} \sqrt{\frac{1}{L_2 \Theta_1}}, \frac{\varepsilon}{\sqrt{n}} \frac{1}{\sqrt{L_2 \Theta_1}}\right\}\right)$	$O\left(\min\left\{\frac{\varepsilon^{3/4}}{\sqrt{n}} \sqrt{\frac{1}{L_2 \Theta_2}}, \frac{\varepsilon}{\sqrt{n}} \frac{1}{\sqrt{L_2 \Theta_2}}\right\}\right)$
Δ	$O\left(\min\left\{\frac{\varepsilon^{3/2}}{(n \ln n)^{1/2}} \sqrt{\frac{L_2}{\Theta_1}}, \frac{\varepsilon^2}{n} \frac{1}{\Theta_1}\right\}\right)$	$O\left(\min\left\{\frac{\varepsilon^{3/2}}{n} \sqrt{\frac{L_2}{\Theta_2}}, \frac{\varepsilon^2}{n} \frac{1}{\Theta_2}\right\}\right)$
Func. eval-s	$O\left(\max\left\{\sqrt{\frac{n \ln n L_2 \Theta_1}{\varepsilon}}, \frac{\sigma^2 \Theta_1 \ln n}{\varepsilon^2}\right\}\right)$	$O\left(\max\left\{\sqrt{\frac{n^2 L_2 \Theta_2}{\varepsilon}}, \frac{\sigma^2 \Theta_2 n}{\varepsilon^2}\right\}\right)$

Table 1: Algorithm 1 parameters for the cases $p = 1$ and $p = 2$.**Algorithm 2** Randomized Derivative-Free Directional Search (RDFDS)

Input: x_0 —starting point; N — number of iterations; $m \geq 1$ — batch size; $t > 0$ — smoothing parameter.

Output: point \bar{x}_N .

- 1: **for** $k = 0, \dots, N - 1$. **do**
- 2: $\alpha \leftarrow \frac{1}{48n\rho_n L_2}$.
- 3: Generate $e_{k+1} \in RS_2(1)$ independently from previous iterations and $\xi_i, i = 1, \dots, m$ — independent realizations of ξ .
- 4: Calculate $\tilde{\nabla}^m f^t(x_{k+1})$ given in (9).
- 5: $x_{k+1} \leftarrow \operatorname{argmin}_{x \in \mathbb{R}^n} \left\{ \alpha n \left\langle \tilde{\nabla}^m f^t(x_k), x - x_k \right\rangle + V[x_k](x) \right\}$.
- 6: **end for**
- 7: **return** $\bar{x}_N \leftarrow \frac{1}{N} \sum_{k=0}^{N-1} x_k$

Theorem 3 Let RDFDS with $t = 2\sqrt{\frac{\Delta}{L_2}}$ be applied to solve problem (1) and $\Theta_p = V[x_0](x^*)$. Then

$$\mathbb{E}[f(\bar{x}_N)] - f(x_*) \leq \frac{384n\rho_n L_2 \Theta_p}{N} + \frac{2\sigma^2}{L_2 m} + \frac{2n\Delta}{3} + \frac{16\sqrt{2n\Theta_p \Delta L_2}}{N} + \frac{4\Delta N}{3\rho_n}, \quad (11)$$

The proof of the theorem is given in the Section 3. As before, we summarize the appropriate values of the algorithm parameters in Table 2 below.

2.3. Inequalities for Gradient Approximation

The proof of the main theorem relies on the following technical result, which connects finite-difference approximation (9) of the stochastic gradient with the stochastic gradient itself and also with ∇f .

	$p = 1$	$p = 2$
N	$O\left(\frac{nL_2\Theta_1}{\varepsilon}\right)$	$O\left(\frac{\ln nL_2\Theta_2}{\varepsilon}\right)$
m	$O\left(\max\left\{1, \frac{\sigma^2}{L_2\varepsilon}\right\}\right)$	$O\left(\max\left\{1, \frac{\sigma^2}{L_2\varepsilon}\right\}\right)$
Δ	$O\left(\min\left\{\frac{\varepsilon}{n}, \frac{\varepsilon^2}{nL_2\Theta_1}\right\}\right)$	$O\left(\min\left\{\frac{\varepsilon}{n}, \frac{\varepsilon^2}{nL_2\Theta_2}\right\}\right)$
Func. eval-s	$O\left(\max\left\{\frac{nL_2\Theta_1}{\varepsilon}, \frac{\sigma^2\Theta_1n}{\varepsilon^2}\right\}\right)$	$O\left(\max\left\{\frac{nL_2\Theta_2}{\varepsilon}, \frac{\sigma^2\Theta_2n}{\varepsilon^2}\right\}\right)$

Table 2: Algorithm 2 parameters for the cases $p = 1$ and $p = 2$.

Lemma 4 For all $x, s \in \mathbb{R}^n$, we have

$$\mathbb{E}_e \|\tilde{\nabla}^m f^t(x)\|_q^2 \leq \frac{12\rho_n}{n} \|g^m(x, \boldsymbol{\xi}_m)\|_2^2 + \frac{\rho_n t^2}{m} \sum_{i=1}^m L(\xi_i)^2 + \frac{16\rho_n \Delta^2}{t^2}, \quad (12)$$

$$\mathbb{E}_e \|\tilde{\nabla}^m f^t(x)\|_2^2 \geq \frac{1}{2n} \|g^m(x, \boldsymbol{\xi}_m)\|_2^2 - \frac{t^2}{2m} \sum_{i=1}^m L(\xi_i)^2 - \frac{8\Delta^2}{t^2}, \quad (13)$$

$$\mathbb{E}_e \langle \tilde{\nabla}^m f^t(x), s \rangle \geq \frac{1}{n} \langle g^m(x, \boldsymbol{\xi}_m), s \rangle - \frac{t\|s\|_p}{2m\sqrt{n}} \sum_{i=1}^m L(\xi_i) - \frac{2\Delta\|s\|_p}{t\sqrt{n}}, \quad (14)$$

$$\mathbb{E}_e \|\langle \nabla f(x), e \rangle e - \tilde{\nabla}^m f^t(x)\|_2^2 \leq \frac{2}{n} \|\nabla f(x) - g^m(x, \boldsymbol{\xi}_m)\|_2^2 + \frac{t^2}{m} \sum_{i=1}^m L(\xi_i)^2 + \frac{16\Delta^2}{t^2}, \quad (15)$$

where $g^m(x, \boldsymbol{\xi}_m) := \frac{1}{m} \sum_{i=1}^m g(x, \xi_i)$, Δ is defined in (3), $L(\xi)$ is the Lipschitz constant of $g(\cdot, \xi)$, which is the gradient of $F(\cdot, \xi)$.

Proof First of all, we rewrite $\tilde{\nabla}^m f^t(x)$ as follows

$$\tilde{\nabla}^m f^t(x) = \left(\langle g^m(x, \boldsymbol{\xi}_m), e \rangle + \frac{1}{m} \sum_{i=1}^m \theta(x, \xi_i, t, e) \right) e,$$

where

$$\theta(x, \xi_i, t, e) = \frac{F(x + te, \xi_i) - F(x, \xi_i)}{t} - \langle g(x, \xi_i), e \rangle + \frac{\Delta(x + te, \xi_i) - \Delta(x, \xi_i)}{t}, \quad i = 1, \dots, m.$$

By Lipschitz smoothness of $F(\cdot, \xi)$ and (3), we have

$$|\theta(x, \xi_i, t, e)| \leq \frac{L(\xi)t}{2} + \frac{2\Delta}{t}. \quad (16)$$

Proof of (12).

$$\begin{aligned}
\mathbb{E}_e \|\tilde{\nabla}^m f^t(x)\|_q^2 &= \mathbb{E}_e \left\| \left(\langle g^m(x, \xi_m), e \rangle + \frac{1}{m} \sum_{i=1}^m \theta(x, \xi_i, t, e) \right) e \right\|_q^2 \\
&\stackrel{\textcircled{1}}{\leq} 2\mathbb{E}_e \|\langle g^m(x, \xi_m), e \rangle e\|_q^2 + 2\mathbb{E}_e \left\| \frac{1}{m} \sum_{i=1}^m \theta(x, \xi_i, t, e) e \right\|_q^2 \\
&\stackrel{\textcircled{2}}{\leq} \frac{12\rho_n}{n} \|g^m(x, \xi_m)\|_2^2 + \frac{2\rho_n}{m} \sum_{i=1}^m \left(\frac{L(\xi_i)t}{2} + \frac{2\Delta}{t} \right)^2 \leq \frac{12\rho_n}{n} \|g^m(x, \xi_m)\|_2^2 + \frac{\rho_n t^2}{m} \sum_{i=1}^m L(\xi_i)^2 + \frac{16\rho_n \Delta^2}{t^2},
\end{aligned} \tag{17}$$

where $\textcircled{1}$ holds since $\|x + y\|_q^2 \leq 2\|x\|_q^2 + 2\|y\|_q^2, \forall x, y \in \mathbb{R}^n$; $\textcircled{2}$ follows from inequalities (7), (8), (16) and the fact that, for any $a_1, a_2, \dots, a_m > 0$, it holds that $\left(\sum_{i=1}^m a_i \right)^2 \leq m \sum_{i=1}^m a_i^2$.

Proof of (13).

$$\begin{aligned}
\mathbb{E}_e \|\tilde{\nabla}^m f^t(x)\|_2^2 &= \mathbb{E}_e \left\| \left(\langle g^m(x, \xi_m), e \rangle + \frac{1}{m} \sum_{i=1}^m \theta(x, \xi_i, t, e) \right) e \right\|_2^2 \\
&\stackrel{\textcircled{1}}{\geq} \frac{1}{2} \mathbb{E}_e \|\langle g^m(x, \xi_m), e \rangle e\|_2^2 - \frac{1}{m} \sum_{i=1}^m \left(\frac{L(\xi_i)t}{2} + \frac{2\Delta}{t} \right)^2 \stackrel{\textcircled{2}}{\geq} \frac{1}{2n} \|g^m(x, \xi_m)\|_2^2 - \frac{t^2}{2m} \sum_{i=1}^m L(\xi_i)^2 - \frac{8\Delta^2}{t^2},
\end{aligned} \tag{18}$$

where $\textcircled{1}$ follows from (16) and inequality $\|x + y\|_2^2 \geq \frac{1}{2}\|x\|_2^2 - \|y\|_2^2, \forall x, y \in \mathbb{R}^n$; $\textcircled{2}$ follows from $e \in S_2(1)$ and Lemma B.10 in [Bogolubsky et al. \(2016\)](#), stating that, for any $s \in \mathbb{R}^n$, $\mathbb{E} \langle s, e \rangle^2 = \frac{1}{n} \|s\|_2^2$.

Proof of (14).

$$\begin{aligned}
\mathbb{E}_e \langle \tilde{\nabla}^m f^t(x), s \rangle &= \mathbb{E}_e \langle \langle g^m(x, \xi_m), e \rangle e, s \rangle + \mathbb{E}_e \frac{1}{m} \sum_{i=1}^m \theta(x, \xi_i, t, e) \langle e, s \rangle \\
&\stackrel{\textcircled{1}}{\geq} \frac{1}{n} \langle g^m(x, \xi_m), s \rangle - \frac{1}{m} \sum_{i=1}^m \left(\frac{L(\xi_i)t}{2} + \frac{2\Delta}{t} \right) \mathbb{E}_e |\langle e, s \rangle| \\
&\stackrel{\textcircled{2}}{\geq} \frac{1}{n} \langle g^m(x, \xi_m), s \rangle - \frac{t\|s\|_p}{2m\sqrt{n}} \sum_{i=1}^m L(\xi_i) - \frac{2\Delta\|s\|_p}{t\sqrt{n}}
\end{aligned} \tag{19}$$

where $\textcircled{1}$ follows from $\mathbb{E}_e [n \langle g, e \rangle e] = g, \forall g \in \mathbb{R}^n$ and (16); $\textcircled{2}$ follows from Lemma B.10 in [Bogolubsky et al. \(2016\)](#), since $\mathbb{E} |\langle s, e \rangle| \leq \sqrt{\mathbb{E} \langle s, e \rangle^2}$, and the fact that $\|x\|_2 \leq \|x\|_p$ for $p \leq 2$.

Proof of (15).

$$\begin{aligned}
\mathbb{E}_e \|\langle \nabla f(x), e \rangle e - \tilde{\nabla}^m f^t(x)\|_2^2 &= \mathbb{E}_e \left\| \langle \nabla f(x), e \rangle e - \langle g^m(x, \xi_m), e \rangle e - \frac{1}{m} \sum_{i=1}^m \theta(x, \xi_i, t, e) e \right\|_2^2 \\
&\stackrel{\textcircled{1}}{\leq} 2\mathbb{E}_e \|\langle \nabla f(x) - g^m(x, \xi_m), e \rangle e\|_2^2 + 2\mathbb{E}_e \left\| \frac{1}{m} \sum_{i=1}^m \theta(x, \xi_i, t, e) e \right\|_2^2 \\
&\stackrel{\textcircled{2}}{\leq} \frac{2}{n} \|\nabla f(x) - g^m(x, \xi_m)\|_2^2 + \frac{t^2}{m} \sum_{i=1}^m L(\xi_i)^2 + \frac{16\Delta^2}{t^2},
\end{aligned} \tag{20}$$

where $\textcircled{1}$ holds since $\|x + y\|_2^2 \leq 2\|x\|_2^2 + 2\|y\|_2^2, \forall x, y \in \mathbb{R}^n$; $\textcircled{2}$ follows from $e \in S_2(1)$ and Lemma B.10 in [Bogolubsky et al. \(2016\)](#), and (16). ■

2.4. Estimates for the Progress of the Method

The following lemma estimates the progress in step 7 of ARDFDS, which is a gradient step.

Lemma 5 Assume that $y = x - \frac{1}{2L_2} \tilde{\nabla}^m f^t(x)$. Then,

$$\begin{aligned} \|g^m(x, \xi_m)\|_2^2 &\leq 8nL_2(f(x) - \mathbb{E}_e f(y)) + 8\|\nabla f(x) - g^m(x, \xi_m)\|_2^2 \\ &\quad + \frac{5nt^2}{m} \sum_{i=1}^m L(\xi_i)^2 + \frac{80n\Delta^2}{t^2}, \end{aligned} \quad (21)$$

where $g^m(x, \xi_m)$ is defined in Lemma 4, Δ is defined in (3), $L(\xi)$ is the Lipschitz constant of $g(\cdot, \xi)$, which is the gradient of $F(\cdot, \xi)$.

Proof Since $\tilde{\nabla}^m f^t(x)$ is collinear to e , we have that, for some $\gamma \in \mathbb{R}$, $y - x = \gamma e$. Then, since $\|e\|_2 = 1$,

$$\langle \nabla f(x), y - x \rangle = \langle \nabla f(x), e \rangle \gamma = \langle \nabla f(x), e \rangle \langle e, y - x \rangle = \langle \nabla f(x), e \rangle e, y - x \rangle.$$

From this and L_2 -smoothness of f we obtain

$$\begin{aligned} f(y) &\leq f(x) + \langle \nabla f(x), e \rangle e, y - x \rangle + \frac{L_2}{2} \|y - x\|_2^2 \\ &\leq f(x) + \langle \tilde{\nabla}^m f^t(x), y - x \rangle + L_2 \|y - x\|_2^2 + \langle \nabla f(x), e \rangle e - \tilde{\nabla}^m f^t(x), y - x \rangle - \frac{L_2}{2} \|y - x\|_2^2 \\ &\stackrel{\textcircled{1}}{\leq} f(x) + \langle \tilde{\nabla}^m f^t(x), y - x \rangle + L_2 \|y - x\|_2^2 + \frac{1}{2L_2} \|\nabla f(x), e \rangle e - \tilde{\nabla}^m f^t(x)\|_2^2, \end{aligned}$$

where $\textcircled{1}$ follows from the Fenchel inequality $\langle s, z \rangle - \frac{\zeta}{2} \|z\|_2^2 \leq \frac{1}{2\zeta} \|s\|_2^2$. Using $y = x - \frac{1}{2L_2} \tilde{\nabla}^m f^t(x)$, we get

$$\frac{1}{4L_2} \|\tilde{\nabla}^m f^t(x)\|_2^2 \leq f(x) - f(y) + \frac{1}{2L_2} \|\nabla f(x), e \rangle e - \tilde{\nabla}^m f^t(x)\|_2^2$$

Taking the expectation in e and applying (13), (15), we obtain

$$\begin{aligned} \frac{1}{4L_2} \left(\frac{1}{2n} \|g^m(x, \xi_m)\|_2^2 - \frac{t^2}{2m} \sum_{i=1}^m L(\xi_i)^2 - \frac{8\Delta^2}{t^2} \right) &\leq \frac{1}{4L_2} \mathbb{E}_e \|\tilde{\nabla}^m f^t(x)\|_2^2 \\ &\leq f(x) - \mathbb{E}_e f(y) + \frac{1}{2L_2} \mathbb{E}_e \|\nabla f(x), e \rangle e - \tilde{\nabla}^m f^t(x)\|_2^2 \\ &\leq f(x) - \mathbb{E}_e f(y) + \frac{1}{2L_2} \left(\frac{2}{n} \|\nabla f(x) - g^m(x, \xi_m)\|_2^2 + \frac{t^2}{m} \sum_{i=1}^m L(\xi_i)^2 + \frac{16\Delta^2}{t^2} \right), \end{aligned}$$

Rearranging the terms, we obtain the statement of the lemma. ■

The following lemma estimates the progress in step 8 of ARDFDS, which is a Mirror Descent step.

Lemma 6 Assume that $z_+ = \operatorname{argmin}_{u \in \mathbb{R}^n} \left\{ \alpha n \left\langle \tilde{\nabla}^m f^t(x), u - z \right\rangle + V[z](u) \right\}$. Then,

$$\begin{aligned} \alpha \langle g^m(x, \xi_m), z - u \rangle &\leq 6\alpha^2 n \rho_n \|g^m(x, \xi_m)\|_2^2 + V[z](u) - \mathbb{E}_e [V[z_+](u)] \\ &\quad + \frac{\alpha^2 n^2 \rho_n}{2} \left(\frac{t^2}{m} \sum_{i=1}^m L(\xi_i)^2 + \frac{16\Delta^2}{t^2} \right) + \alpha \sqrt{n} \|z - u\|_p \left(\frac{t}{2m} \sum_{i=1}^m L(\xi_i) + \frac{2\Delta}{t} \right), \end{aligned} \quad (22)$$

where $g^m(x, \xi_m)$ is defined in Lemma 4, Δ is defined in (3), $L(\xi)$ is the Lipschitz constant of $g(\cdot, \xi)$, which is the gradient of $F(\cdot, \xi)$.

Proof For all $u \in \mathbb{R}^n$, we have

$$\begin{aligned}
\alpha n \langle \tilde{\nabla}^m f^t(x), z - u \rangle &= \alpha n \langle \tilde{\nabla}^m f^t(x), z - z_+ \rangle + \alpha n \langle \tilde{\nabla}^m f^t(x), z_+ - u \rangle \\
&\stackrel{\textcircled{1}}{\leq} \alpha n \langle \tilde{\nabla}^m f^t(x), z - z_+ \rangle + \langle -\nabla V[z](z_+), z_+ - u \rangle \stackrel{\textcircled{2}}{=} \alpha n \langle \tilde{\nabla}^m f^t(x), z - z_+ \rangle \\
&+ V[z](u) - V[z_+](u) - V[z](z_+) \stackrel{\textcircled{3}}{\leq} \left(\alpha n \langle \tilde{\nabla}^m f^t(x), z - z_+ \rangle - \frac{1}{2} \|z - z_+\|_p^2 \right) \\
&+ V[z](u) - V[z_+](u) \stackrel{\textcircled{4}}{\leq} \frac{\alpha^2 n^2}{2} \|\tilde{\nabla}^m f^t(x)\|_q^2 + V[z](u) - V[z_+](u),
\end{aligned} \tag{23}$$

where $\textcircled{1}$ follows from the definition of z_+ , whence $\langle \nabla V[z](z_+) + \alpha n \tilde{\nabla}^m f^t(x), u - z_+ \rangle \geq 0$ for all $u \in \mathbb{R}^n$; $\textcircled{2}$ follows from the "magic identity" Fact 5.3.3 in Ben-Tal and Nemirovski (2015) for the Bregman divergence; $\textcircled{3}$ follows from (6); and $\textcircled{4}$ follows from the Fenchel inequality $\zeta \langle s, z \rangle - \frac{1}{2} \|z\|_p^2 \leq \frac{\zeta^2}{2} \|s\|_q^2$. Taking expectation in e and applying (14) with $s = z - u$, (12), we get

$$\begin{aligned}
\alpha n \left(\frac{1}{n} \langle g^m(x, \xi_m), z - u \rangle - \frac{t \|z - u\|_p}{2m\sqrt{n}} \sum_{i=1}^m L(\xi_i) - \frac{2\Delta \|z - u\|_p}{t\sqrt{n}} \right) &\leq \alpha n \mathbb{E}_e \langle \tilde{\nabla}^m f^t(x), z - u \rangle \\
&\leq \frac{\alpha^2 n^2}{2} \mathbb{E}_e \|\tilde{\nabla}^m f^t(x)\|_q^2 + V[z](u) - \mathbb{E}_e[V[z_+](u)] \\
&\leq \frac{\alpha^2 n^2}{2} \left(\frac{12\rho_n}{n} \|g^m(x, \xi_m)\|_2^2 + \frac{\rho_n t^2}{m} \sum_{i=1}^m L(\xi_i)^2 + \frac{16\rho_n \Delta^2}{t^2} \right) + V[z](u) - \mathbb{E}_e[V[z_+](u)].
\end{aligned} \tag{24}$$

Rearranging the terms, we obtain the statement of the lemma. \blacksquare

2.5. Proof of Theorem 2

First, we prove the following lemma, which estimates the one-iteration progress of the whole algorithm.

Lemma 7 Let $\{x_k, y_k, z_k, \alpha_k, \tau_k\}$, $k \geq 0$ be generated by ARDFDS. Then, for all $u \in \mathbb{R}^n$,

$$\begin{aligned}
48n^2 \rho_n L_2 \alpha_{k+1}^2 \mathbb{E}_{e, \xi} [f(y_{k+1}) \mid \mathcal{E}_k, \Xi_k] - (48n^2 \rho_n L_2 \alpha_{k+1}^2 - \alpha_{k+1}) f(y_k) \\
- V[z_k](u) + \mathbb{E}_{e, \xi} [V[z_{k+1}](u) \mid \mathcal{E}_k, \Xi_k] - \mathcal{R}_{k+1} \leq \alpha_{k+1} f(u),
\end{aligned} \tag{25}$$

$$\mathcal{R}_{k+1} := 48\alpha_{k+1}^2 n \rho_n \frac{\sigma^2}{m} + \frac{61\alpha_{k+1}^2 n^2 \rho_n}{2} \left(L_2^2 t^2 + \frac{16\Delta^2}{t^2} \right) + \alpha_{k+1} \sqrt{n} \|z_k - u\|_p \left(\frac{L_2 t}{2} + \frac{2\Delta}{t} \right). \tag{26}$$

where Δ is defined in (3), \mathcal{E}_k and Ξ_k denote the history of realizations of e_1, \dots, e_k and $\xi_{1,1}, \dots, \xi_{k,m}$ respectively, up to the step k .

Proof Combining (21) and (22), we obtain

$$\begin{aligned}
\alpha \langle g^m(x_{k+1}, \xi_m^{(k+1)}), z - u \rangle &\leq 48\alpha^2 n^2 \rho_n L_2 (f(x_{k+1}) - \mathbb{E}_e f(y_{k+1})) \\
&+ V[z_k](u) - \mathbb{E}_e[V[z_{k+1}](u)] + 48\alpha^2 n \rho_n \|\nabla f(x_{k+1}) - g^m(x_{k+1}, \xi_m^{(k+1)})\|_2^2 \\
&+ \frac{61\alpha^2 n^2 \rho_n}{2} \left(\frac{t^2}{m} \sum_{i=1}^m L(\xi_i^{(k+1)})^2 + \frac{16\Delta^2}{t^2} \right) + \alpha \sqrt{n} \|z_k - u\|_p \left(\frac{t}{2m} \sum_{i=1}^m L(\xi_i^{(k+1)}) + \frac{2\Delta}{t} \right),
\end{aligned} \tag{27}$$

where $g^m(x, \xi_m)$ is defined in Lemma 4 and the expectation in e is conditional to \mathcal{E}_k . By the definition of $g^m(x, \xi_m)$ and (2), $\mathbb{E}_\xi g^m(x, \xi_m) = \nabla f(x)$ and $\mathbb{E}_\xi \|\nabla f(x_{k+1}) - g^m(x_{k+1}, \xi_m^{(k+1)})\|_2^2 \leq \frac{\sigma^2}{m}$. Using these two facts and taking the expectation in $\xi_m^{(k+1)}$ conditional to Ξ_k , we obtain

$$\begin{aligned}
\alpha_{k+1} \langle \nabla f(x_{k+1}), z_k - u \rangle &\leq 48\alpha_{k+1}^2 n^2 \rho_n L_2 (f(x_{k+1}) - \mathbb{E}_{e, \xi} [f(y_{k+1}) \mid \mathcal{E}_k, \Xi_k]) \\
&+ V[z_k](u) - \mathbb{E}_{e, \xi} [V[z_{k+1}](u) \mid \mathcal{E}_k, \Xi_k] + \mathcal{R}_{k+1}.
\end{aligned} \tag{28}$$

Further,

$$\begin{aligned}
& \alpha_{k+1}(f(x_{k+1}) - f(u)) \leq \alpha_{k+1}\langle \nabla f(x_{k+1}), x_{k+1} - u \rangle \\
& = \alpha_{k+1}\langle \nabla f(x_{k+1}), x_{k+1} - z_k \rangle + \alpha_{k+1}\langle \nabla f(x_{k+1}), z_k - u \rangle \\
& \stackrel{\textcircled{1}}{=} \frac{(1-\tau_k)\alpha_{k+1}}{\tau_k} \langle \nabla f(x_{k+1}), y_k - x_{k+1} \rangle + \alpha_{k+1}\langle \nabla f(x_{k+1}), z_k - u \rangle \\
& \stackrel{\textcircled{2}}{\leq} \frac{(1-\tau_k)\alpha_{k+1}}{\tau_k} (f(y_k) - f(x_{k+1})) + \alpha_{k+1}\langle \nabla f(x_{k+1}), z_k - u \rangle \\
& \stackrel{\textcircled{28}}{\leq} \frac{(1-\tau_k)\alpha_{k+1}}{\tau_k} (f(y_k) - f(x_{k+1})) + 48\alpha_{k+1}^2 n^2 \rho_n L_2 (f(x_{k+1}) - \mathbb{E}_{e_{k+1}, \xi_{k+1}}[f(y_{k+1}) \mid \mathcal{E}_k, \Xi_k]) \\
& \quad + V[z_k](u) - \mathbb{E}_{e_{k+1}, \xi_{k+1}}[V[z_{k+1}](u) \mid \mathcal{E}_k, \Xi_k] + \mathcal{R}_{k+1} \\
& \stackrel{\textcircled{3}}{=} (48\alpha_{k+1}^2 n^2 \rho_n L_2 - \alpha_{k+1})f(y_k) - 48\alpha_{k+1}^2 n^2 \rho_n L_2 \mathbb{E}_{e_{k+1}}[f(y_{k+1}) \mid \mathcal{E}_k, \Xi_k] \\
& \quad + \alpha_{k+1}f(x_{k+1}) + V[z_k](u) - \mathbb{E}_{e_{k+1}, \xi_{k+1}}[V[z_{k+1}](u) \mid \mathcal{E}_k, \Xi_k] + \mathcal{R}_{k+1}.
\end{aligned}$$

That is, ① is since $x_{k+1} := \tau_k z_k + (1-\tau_k)y_k \Leftrightarrow \tau_k(x_{k+1} - z_k) = (1-\tau_k)(y_k - x_{k+1})$, ② follows from the convexity of f and inequality $1 - \tau_k \geq 0$ and ③ is since $\tau_k = \frac{1}{48\alpha_{k+1}n^2\rho_n L_2}$. Rearranging the terms, we obtain the statement of the lemma. \blacksquare

Proof [Proof of Theorem 2.]

Note that $48n^2\rho_n L_2\alpha_{k+1}^2 - \alpha_{k+1} + \frac{1}{192n^2\rho_n L_2} = 48n^2\rho_n L_2\alpha_k^2$. That is,

$$\begin{aligned}
48n^2\rho_n L_2\alpha_{k+1}^2 - \alpha_{k+1} + \frac{1}{192n^2\rho_n L_2} &= \frac{(k+2)^2}{192n^2\rho_n L_2} - \frac{k+2}{96n^2\rho_n L_2} + \frac{1}{192n^2\rho_n L_2} \\
&= \frac{k^2+4k+4-2k-4+1}{192n^2\rho_n L_2} = \frac{(k+1)^2}{192n^2\rho_n L_2} = 48n^2\rho_n L_2\alpha_k^2.
\end{aligned}$$

Taking full expectation $\mathbb{E}[\cdot] = \mathbb{E}_{e_1, \dots, e_N, \xi_{1,1}, \dots, \xi_{N,m}}[\cdot]$ from both sides of (25) for $k = 0, \dots, l-1$, where $l \leq N$ and telescoping the obtained inequalities⁴ we have

$$\begin{aligned}
& 48n^2\rho_n L_2\alpha_l^2 \mathbb{E}[f(y_l)] + \sum_{k=1}^{l-1} \frac{1}{192n^2\rho_n L_2} \mathbb{E}[f(y_k)] - V[z_0](u) + \mathbb{E}[V[z_l](u)] \\
& - \zeta_1 \sum_{k=0}^{l-1} \alpha_{k+1} \mathbb{E}[\|u - z_k\|_p] - \zeta_2 \sum_{k=0}^{l-1} \alpha_{k+1}^2 \leq \sum_{k=0}^{l-1} \alpha_{k+1} f(u),
\end{aligned} \tag{29}$$

where we denoted

$$\zeta_1 := \sqrt{n} \left(\frac{L_2 t}{2} + \frac{2\Delta}{t} \right), \quad \zeta_2 := 48n\rho_n \frac{\sigma^2}{m} + \frac{61n^2\rho_n}{2} \left(L_2^2 t^2 + \frac{16\Delta^2}{t^2} \right). \tag{30}$$

We set $u = x^*$ in (29), where x^* is a solution to (1), and define $\Theta := V[z_0](x^*)$, $R_k := \mathbb{E}[\|x^* - z_k\|_p]$. Also, from (6), we have that $\zeta_1 \alpha_1 R_0 \leq \frac{\sqrt{2}\Theta\zeta_1}{48n^2\rho_n L_2}$. To simplify the notation, we define $B_l := \zeta_2 \sum_{k=0}^{l-1} \alpha_{k+1}^2 + \Theta + \frac{\sqrt{2}\Theta\zeta_1}{48n^2\rho_n L_2}$. Since $\sum_{k=0}^{l-1} \alpha_{k+1} = \frac{l(l+3)}{192n^2\rho_n L_2}$ and, for all $i = 1, \dots, N$, $f(y_i) \leq f(x^*)$, we obtain from (29)

$$\begin{aligned}
\frac{(l+1)^2}{192n^2\rho_n L_2} \mathbb{E}[f(y_l)] &\leq f(x^*) \left(\frac{(l+3)l}{192n^2\rho_n L_2} - \frac{l-1}{192n^2\rho_n L_2} \right) + B_l - \mathbb{E}[V[z_l](x^*)] + \zeta_1 \sum_{k=1}^{l-1} \alpha_{k+1} R_k, \\
0 &\leq \frac{(l+1)^2}{192n^2\rho_n L_2} (\mathbb{E}[f(y_l)] - f(x^*)) \leq B_l - \mathbb{E}[V[z_l](x^*)] + \zeta_1 \sum_{k=1}^{l-1} \alpha_{k+1} R_k,
\end{aligned} \tag{31}$$

4. Note that $\alpha_1 = \frac{2}{96n^2\rho_n L_2} = \frac{1}{48n^2\rho_n L_2}$ and therefore $48n^2\rho_n L_2\alpha_1^2 - \alpha_1 = 0$.

which gives

$$\mathbb{E}[V[z_l](x^*)] \leq B_l + \zeta_1 \sum_{k=1}^{l-1} \alpha_{k+1} R_k. \quad (32)$$

Moreover,

$$\frac{1}{2} (\mathbb{E}[\|z_l - x^*\|_p])^2 \leq \frac{1}{2} \mathbb{E}[\|z_l - x^*\|_p^2] \leq \mathbb{E}[V[z_l](x^*)] \stackrel{(32)}{\leq} B_l + \zeta_1 \sum_{k=1}^{l-1} \alpha_{k+1} R_k, \quad (33)$$

whence,

$$R_l \leq \sqrt{2} \cdot \sqrt{B_l + \zeta_1 \sum_{k=1}^{l-1} \alpha_{k+1} R_k}. \quad (34)$$

Applying Lemma 8 for $a_0 = \zeta_2 \alpha_1^2 + \Theta + \frac{\sqrt{2\Theta}\zeta_1}{48n^2\rho_n L_2}$, $a_k = \zeta_2 \alpha_{k+1}^2$, $b = \zeta_1$ for $k = 1, \dots, N-1$, we obtain

$$B_l + \zeta_1 \sum_{k=1}^{l-1} \alpha_{k+1} R_k \leq \left(\sqrt{B_l} + \sqrt{2}\zeta_1 \cdot \frac{l^2}{96n^2\rho_n L_2} \right)^2, \quad l = 1, \dots, N \quad (35)$$

Since $V[z](x^*) \geq 0$, by inequality (31) for $l = N$ and the definition of B_l , we have

$$\begin{aligned} \frac{(N+1)^2}{192n^2\rho_n L_2} (\mathbb{E}[f(y_N)] - f(x^*)) &\leq \left(\sqrt{B_N} + \sqrt{2}\zeta_1 \cdot \frac{N^2}{96n^2\rho_n L_2} \right)^2 \stackrel{\textcircled{1}}{\leq} 2B_N + 4\zeta_1^2 \cdot \frac{N^4}{(96n^2\rho_n L_2)^2} \\ &= 2\zeta_2 \sum_{k=0}^{l-1} \alpha_{k+1}^2 + 2\Theta + \frac{\sqrt{2\Theta}\zeta_1}{24n^2\rho_n L_2} + 4\zeta_1^2 \cdot \frac{N^4}{(96n^2\rho_n L_2)^2} \\ &\stackrel{\textcircled{2}}{\leq} 2\Theta + \frac{\sqrt{2\Theta}\zeta_1}{24n^2\rho_n L_2} + \frac{2\zeta_2(N+1)^3}{(96n^2\rho_n L_2)^2} + 4\zeta_1^2 \cdot \frac{N^4}{(96n^2\rho_n L_2)^2} \end{aligned} \quad (36)$$

where $\textcircled{1}$ is due to the fact that $\forall a, b \in \mathbb{R} \quad (a+b)^2 \leq 2a^2 + 2b^2$ and $\textcircled{2}$ is because $\sum_{k=0}^{N-1} \alpha_{k+1}^2 =$

$$\frac{1}{(96n^2\rho_n L_2)^2} \sum_{k=2}^{N+1} k^2 \leq \frac{1}{(96n^2\rho_n L_2)^2} \cdot \frac{(N+1)(N+2)(2N+3)}{6} \leq \frac{1}{(96n^2\rho_n L_2)^2} \cdot \frac{(N+1)2(N+1)3(N+1)}{6} = \frac{(N+1)^3}{(96n^2\rho_n L_2)^2}.$$

Dividing (36) by $\frac{(N+1)^2}{192n^2\rho_n L_2}$ and substituting ζ_1, ζ_2 from (30), we obtain

$$\begin{aligned} \mathbb{E}[f(y_N)] - f(x^*) &\leq \frac{384\Theta n^2\rho_n L_2}{(N+1)^2} + \frac{12\sqrt{2\Theta}}{(N+1)^2} \zeta_1 + \frac{384(N+1)\zeta_2}{(96n^2\rho_n L_2)^2} + \frac{N^4\zeta_1^2}{12n^2\rho_n L_2(N+1)^2} \\ &\leq \frac{384\Theta n^2\rho_n L_2}{N^2} + \frac{12\sqrt{2n\Theta}}{N^2} \left(\frac{L_2 t}{2} + \frac{2\Delta}{t} \right) + \frac{384N}{nL_2} \frac{\sigma^2}{m} \\ &\quad + \frac{6N}{L_2} \left(L_2^2 t^2 + \frac{16\Delta^2}{t^2} \right) + \frac{N^2}{24n\rho_n L_2} \left(L_2^2 t^2 + \frac{16\Delta^2}{t^2} \right). \end{aligned}$$

■

3. Randomized Derivative-Free Directional Search

In this section we prove the convergence rate theorem for Randomized Derivative-Free Directional Search algorithm.

Proof [Proof of Theorem 3]

$$\begin{aligned}
\alpha n \langle \tilde{\nabla}^m f^t(x_k), x_k - x_* \rangle &= \alpha n \langle \tilde{\nabla}^m f^t(x_k), x_k - x_{k+1} \rangle + \alpha n \langle \tilde{\nabla}^m f^t(x_k), x_{k+1} - x_* \rangle \\
&\stackrel{\textcircled{1}}{\leq} \alpha n \langle \tilde{\nabla}^m f^t(x_k), x_k - x_{k+1} \rangle + \langle -\nabla V[x_k](x_{k+1}), x_{k+1} - x_* \rangle \\
&\stackrel{\textcircled{2}}{=} \alpha n \langle \tilde{\nabla}^m f^t(x_k), x_k - x_{k+1} \rangle + V[x_k](x_*) - V[x_{k+1}](x_*) - V[x_k](x_{k+1}) \\
&\stackrel{\textcircled{3}}{\leq} \left(\alpha n \langle \tilde{\nabla}^m f^t(x_k), x_k - x_{k+1} \rangle - \frac{1}{2} \|x_k - x_{k+1}\|_p^2 \right) \\
&\quad + V[x_k](x_*) - V[x_{k+1}](x_*) \leq \frac{\alpha^2 n^2}{2} \|\tilde{\nabla}^m f^t(x_k)\|_q^2 + V[x_k](x_*) - V[x_{k+1}](x_*),
\end{aligned} \tag{37}$$

where $\textcircled{1}$ follows from $x_{k+1} = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \left\{ V[x_k](x) + \alpha \langle n \tilde{\nabla}^m f^t(x_k), x \rangle \right\}$, whence $\langle \nabla V[x_k](x_{k+1}), x_{k+1} - x_* \rangle + \alpha n \langle \tilde{\nabla}^m f^t(x_k), x_{k+1} - x_* \rangle \geq 0$ for all $x \in \mathbb{R}^n$, $\textcircled{2}$ follows from triangle equality for Bregman divergence and $\textcircled{3}$ is due to $V[x](y) \geq \frac{1}{2} \|x - y\|_p^2$. Taking conditional mathematical expectation $\mathbb{E}_{e_{k+1}}[\cdot \mid \mathcal{E}_k]$ from both sides of (37) we get

$$\begin{aligned}
&\alpha n \mathbb{E}_{e_{k+1}}[\langle \tilde{\nabla}^m f^t(x_k), x_k - x_* \rangle \mid \mathcal{E}_k] \\
&\leq \frac{\alpha^2 n^2}{2} \mathbb{E}_{e_{k+1}}[\|\tilde{\nabla}^m f^t(x_k)\|_q^2 \mid \mathcal{E}_k] + V[x_k](x_*) - \mathbb{E}_{e_{k+1}}[V[x_{k+1}](x_*) \mid \mathcal{E}_k]
\end{aligned} \tag{38}$$

From (38), (12) and (14) for $s = x_k - x_*$ we obtain

$$\begin{aligned}
&\langle g^m(x_k, \xi_m^{(k+1)}), x_k - x_* \rangle \leq 24\alpha^2 n \rho_n L_2 (f(x_k) - f(x_*)) \\
&+ 12\alpha^2 n \rho_n \|\nabla f(x_k) - g^m(x_k, \xi_m^{(k+1)})\|_2^2 + \alpha^2 n^2 \rho_n \cdot \frac{t^2}{2m} \sum_{i=1}^m L_2(\xi_{k+1,i})^2 + \frac{8\alpha^2 n^2 \rho_n \Delta^2}{t^2} \\
&\quad + \alpha \sqrt{n} \|x_k - x_*\|_p \cdot \frac{t}{2m} \sum_{i=1}^m L_2(\xi_{k+1,i}) + \frac{2\alpha \Delta \sqrt{n} \|x_k - x_*\|_p}{t} \\
&\quad + V[x_k](x_*) - \mathbb{E}_{e_{k+1}}[V[x_{k+1}](x_*) \mid \mathcal{E}_k].
\end{aligned}$$

Taking conditional mathematical expectation $\mathbb{E}_{\xi_{k+1}}[\cdot \mid \Xi_k] = \mathbb{E}_{\xi_{k+1,1}, \xi_{k+1,2}, \dots, \xi_{k+1,m}}[\cdot \mid \xi_{1,1}, \xi_{1,2}, \dots, \xi_{k,m}]$ from the both sides of previous inequality and using convexity of f and (2) we have

$$\begin{aligned}
&\underbrace{(\alpha - 24\alpha^2 n \rho_n L_2)}_{\frac{\alpha}{4}} (f(x_k) - f(x_*)) \leq 12\alpha^2 n \rho_n \frac{\sigma^2}{m} \\
&+ \alpha^2 n^2 \rho_n \left(\frac{L_2^2 t^2}{2} + \frac{8\Delta^2}{t^2} \right) + \alpha \sqrt{n} \|x_k - x_*\|_p \left(\frac{L_2 t}{2} + \frac{2\Delta}{t} \right) \\
&\quad + V[x_k](x_*) - \mathbb{E}_{e_{k+1}, \xi_{k+1}}[V[x_{k+1}](x_*) \mid \mathcal{E}_k, \Xi_k],
\end{aligned} \tag{39}$$

because $\alpha = \frac{1}{48n\rho_n L_2}$. Denote

$$\zeta_1 = \frac{L_2 t}{2} + \frac{2\Delta}{t}, \quad \zeta_2 = \frac{L_2^2 t^2}{2} + \frac{8\Delta^2}{t^2}. \tag{40}$$

Note that

$$\zeta_1^2 = \left(\frac{L_2 t}{2} + \frac{2\Delta}{t} \right)^2 \leq 2 \cdot \frac{L_2^2 t^2}{4} + 2 \cdot \frac{4\Delta^2}{t^2} = \zeta_2. \tag{41}$$

Taking mathematical expectation $\mathbb{E}[\cdot] = \mathbb{E}_{e_1, e_2, \dots, e_N, \xi_{1,1}, \xi_{1,2}, \dots, \xi_{N,m}}[\cdot]$ from inequalities (39) for $k = 0, \dots, l-1$, where $l \leq N$, and summing them we get

$$\begin{aligned}
0 &\leq \frac{N\alpha}{4} (\mathbb{E}[f(\bar{x}_l)] - f(x_*)) \leq l \cdot 12\alpha^2 n \rho_n \frac{\sigma^2}{m} + l\alpha^2 n^2 \rho_n \zeta_2 \\
&\quad + \alpha \sqrt{n} \zeta_1 \sum_{k=0}^{l-1} \mathbb{E}[\|x_k - x_*\|_p] + \underbrace{V[x_0](x_*) - \mathbb{E}[V[x_l](x_*)]}_{\Theta_p},
\end{aligned} \tag{42}$$

where $\bar{x}_l \stackrel{\text{def}}{=} \frac{1}{l} \sum_{k=0}^{l-1} x_k$. From the previous inequality we get

$$\begin{aligned} \frac{1}{2} (\mathbb{E}[\|x_l - x_*\|_p])^2 &\leq \frac{1}{2} \mathbb{E}[\|x_l - x_*\|_p^2] \leq \mathbb{E}[V[x_l](x_*)] \\ &\leq \Theta_p + l \cdot 12\alpha^2 n \rho_n \frac{\sigma^2}{m} + l\alpha^2 n^2 \rho_n \zeta_2 + \alpha \sqrt{n} \delta \zeta_1 \sum_{k=0}^{l-1} \mathbb{E}[\|x_k - x_*\|_p], \end{aligned} \quad (43)$$

whence $\forall l \leq N$ we obtain

$$\mathbb{E}[\|x_k - x_*\|_p] \leq \sqrt{2} \sqrt{\Theta_p + l \cdot 12\alpha^2 n \rho_n \frac{\sigma^2}{m} + l\alpha^2 n^2 \rho_n \zeta_2 + \alpha \sqrt{n} \zeta_1 \sum_{k=0}^{l-1} \mathbb{E}[\|x_k - x_*\|_p]}. \quad (44)$$

Denote $R_k = \mathbb{E}[\|x^* - x_k\|_p]$ for $k = 0, \dots, N$. Applying Lemma 9 for $a_0 = \Theta_p + \alpha \sqrt{n} \zeta_1 \mathbb{E}[\|x_0 - x_*\|_p] \leq \Theta_p + \alpha \sqrt{2n\Theta_p \zeta_1}$, $a_k = 12\alpha^2 n \rho_n \frac{\sigma^2}{m} + \alpha^2 n^2 \rho_n \zeta_2$, $b = \sqrt{n} \zeta_1$ for $k = 1, \dots, N-1$ we have for $l = N$

$$\begin{aligned} &\frac{N\alpha}{4} (\mathbb{E}[f(\bar{x}_N)] - f(x_*)) \\ &\leq \left(\sqrt{\Theta_p + N \cdot 12\alpha^2 n \rho_n \frac{\sigma^2}{m} + N\alpha^2 n^2 \rho_n \zeta_2 + \alpha \sqrt{2n\Theta_p \zeta_1} + \sqrt{2n} \zeta_1 \alpha N} \right)^2 \\ &\stackrel{\textcircled{1}}{\leq} 2\Theta_p + 24N\alpha^2 n \rho_n \frac{\sigma^2}{m} + 2N\alpha^2 n^2 \rho_n \zeta_2 + 2\alpha \sqrt{2n\Theta_p \zeta_1} + 4n\zeta_1^2 \alpha^2 N^2, \end{aligned}$$

whence

$$\mathbb{E}[f(\bar{x}_N)] - f(x_*) \stackrel{(41)}{\leq} \frac{384n\rho_n L_2 \Theta_p}{N} + \frac{2\sigma^2}{L_2 m} + \frac{n\zeta_2}{6L_2} + \frac{8\sqrt{2n\Theta_p \zeta_1}}{N} + \frac{\zeta_2 N}{3L_2 \rho_n},$$

because $\alpha = \frac{1}{48n\rho_n L_2}$. Remind that $\zeta_2 \stackrel{(40)}{=} \frac{L_2^2 t^2}{2} + \frac{8\Delta^2}{t^2}$. In order to minimize ζ_1 one can choose $t = 2\sqrt{\frac{\Delta}{L_2}}$ and get $\zeta_2 = 4\Delta L_2$. In this setting $\zeta_1 \stackrel{(40)}{=} \frac{L_2 t}{2} + \frac{2\Delta}{t} = 2\sqrt{\Delta L_2}$. Finally we get

$$\mathbb{E}[f(\bar{x}_N)] - f(x_*) \leq \frac{384n\rho_n L_2 \Theta_p}{N} + \frac{2\sigma^2}{L_2 m} + \frac{2n\Delta}{3} + \frac{16\sqrt{2n\Theta_p \Delta L_2}}{N} + \frac{4\Delta N}{3\rho_n}.$$

■

4. Conclusion

In this paper, we propose two new algorithms for smooth stochastic derivative-free optimization with two-point feedback. Our first algorithm is an accelerated one and the second one is a non-accelerated one. Interestingly, despite the traditional choice of 2-norm proximal setup for unconstrained optimization problems, our analysis shows that the method with 1-norm proximal setup has better complexity.

References

Alekh Agarwal, Ofer Dekel, and Lin Xiao. Optimal algorithms for online convex optimization with multi-point bandit feedback. In *COLT 2010 - The 23rd Conference on Learning Theory*, 2010.

- Alekh Agarwal, Dean P Foster, Daniel J Hsu, Sham M Kakade, and Alexander Rakhlin. Stochastic convex optimization with bandit feedback. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 1035–1043. Curran Associates, Inc., 2011. URL <http://papers.nips.cc/paper/4475-stochastic-convex-optimization-with-bandit-feedback>.
- Zeyuan Allen-Zhu. Katyusha: The first direct acceleration of stochastic gradient methods. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2017, pages 1200–1205, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-4528-6. doi: 10.1145/3055399.3055448. URL <http://doi.acm.org/10.1145/3055399.3055448>. arXiv:1603.05953.
- Francis Bach and Vianney Perchet. Highly-smooth zero-th order online optimization. In Vitaly Feldman, Alexander Rakhlin, and Ohad Shamir, editors, *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pages 257–283, Columbia University, New York, New York, USA, 23–26 Jun 2016. PMLR. URL <http://proceedings.mlr.press/v49/bach16.html>.
- Anastasia Bayandina, Alexander Gasnikov, and Anastasia Lagunovskaya. Gradient-free two-points optimal method for non smooth stochastic convex optimization problem with additional small noise. *Automation and remote control*, 79(7), 2018. arXiv:1701.03821.
- Alexandre Belloni, Tengyuan Liang, Hariharan Narayanan, and Alexander Rakhlin. Escaping the local minima via simulated annealing: Optimization of approximately convex functions. In Peter Grnwald, Elad Hazan, and Satyen Kale, editors, *Proceedings of The 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine Learning Research*, pages 240–265, Paris, France, 03–06 Jul 2015. PMLR. URL <http://proceedings.mlr.press/v40/Belloni15.html>.
- Aaron Ben-Tal and Arkadi Nemirovski. *Lectures on Modern Convex Optimization (Lecture Notes)*. Personal web-page of A. Nemirovski, 2015. URL http://www2.isye.gatech.edu/~nemirovs/Lect_ModConvOpt.pdf.
- Lev Bogolubsky, Pavel Dvurechensky, Alexander Gasnikov, Gleb Gusev, Yurii Nesterov, Andrei M Raigorodskii, Aleksey Tikhonov, and Maksim Zhukovskii. Learning supervised pagerank with gradient-based and gradient-free optimization methods. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 4914–4922. Curran Associates, Inc., 2016. URL <http://papers.nips.cc/paper/6565-learning-supervised-pagerank-with-gradient-based-optimization>. arXiv:1603.00717.
- R.P. Brent. *Algorithms for Minimization Without Derivatives*. Dover Books on Mathematics. Dover Publications, 1973. ISBN 9780486419985. URL <https://books.google.de/books?id=6Ay2biHG-GEC>.
- Sébastien Bubeck and Nicoló Cesa-Bianchi. Regret analysis of stochastic and non-stochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012. ISSN 1935-8237. doi: 10.1561/22000000024. URL <http://dx.doi.org/10.1561/22000000024>.

- Sébastien Bubeck, Yin Tat Lee, and Ronen Eldan. Kernel-based methods for bandit convex optimization. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2017*, pages 72–85, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-4528-6. doi: 10.1145/3055399.3055403. URL <http://doi.acm.org/10.1145/3055399.3055403>. arXiv:1607.03084.
- Nicolò Cesa-bianchi, Alex Conconi, and Claudio Gentile. On the generalization ability of on-line learning algorithms. In T. G. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14*, pages 359–366. MIT Press, 2002. URL <http://papers.nips.cc/paper/2113-on-the-generalization-ability-of-on-line-learning>.
- Ofar Dekel, Ronen Eldan, and Tomer Koren. Bandit smooth convex optimization: Improving the bias-variance tradeoff. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 2926–2934. Curran Associates, Inc., 2015. URL <http://papers.nips.cc/paper/5842-bandit-smooth-convex-optimization-improving-tradeoff>.
- Olivier Devolder. Stochastic first order methods in smooth convex optimization. *CORE Discussion Paper 2011/70*, 2011.
- John C. Duchi, Michael I. Jordan, Martin J. Wainwright, and Andre Wibisono. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Trans. Information Theory*, 61(5):2788–2806, 2015. doi: 10.1109/TIT.2015.2409256. URL <https://doi.org/10.1109/TIT.2015.2409256>. arXiv:1312.2139.
- Pavel Dvurechensky and Alexander Gasnikov. Stochastic intermediate gradient method for convex problems with stochastic inexact oracle. *Journal of Optimization Theory and Applications*, 171(1):121–145, 2016. ISSN 1573-2878. doi: 10.1007/s10957-016-0999-6. URL <http://dx.doi.org/10.1007/s10957-016-0999-6>.
- Pavel Dvurechensky, Alexander Gasnikov, and Alexander Tiurin. Randomized similar triangles method: A unifying framework for accelerated randomized optimization methods (coordinate descent, directional search, derivative-free method). *arXiv:1707.08486*, 2017.
- Abraham D. Flaxman, Adam Tauman Kalai, and H. Brendan McMahan. Online convex optimization in the bandit setting: Gradient descent without a gradient. In *Proceedings of the Sixteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '05*, pages 385–394, Philadelphia, PA, USA, 2005. Society for Industrial and Applied Mathematics. ISBN 0-89871-585-7. URL <http://dl.acm.org/citation.cfm?id=1070432.1070486>.
- A. V. Gasnikov, A. A. Lagunovskaya, I. N. Usmanova, and F. A. Fedorenko. Gradient-free proximal methods with inexact oracle for convex stochastic nonsmooth optimization problems on the simplex. *Automation and Remote Control*, 77(11):2018–2034, Nov 2016a. ISSN 1608-3032. doi: 10.1134/S0005117916110114. URL <http://dx.doi.org/10.1134/S0005117916110114>. arXiv:1412.3890.
- A. V. Gasnikov, E. A. Krymova, A. A. Lagunovskaya, I. N. Usmanova, and F. A. Fedorenko. Stochastic online optimization. single-point and multi-point non-linear multi-armed bandits. convex and strongly-convex case. *Automation and Remote Control*, 78

- (2):224–234, Feb 2017. ISSN 1608-3032. doi: 10.1134/S0005117917020035. URL <http://dx.doi.org/10.1134/S0005117917020035>. arXiv:1509.01679.
- Alexander Gasnikov, Pavel Dvurechensky, and Yurii Nesterov. Stochastic gradient methods with inexact oracle. *Proceedings of Moscow Institute of Physics and Technology*, 8(1):41–91, 2016b. In Russian, first appeared in arXiv:1411.4218.
- Saeed Ghadimi and Guanghui Lan. Stochastic first- and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013. doi: 10.1137/120880811. URL <https://doi.org/10.1137/120880811>. arXiv:1309.5549.
- Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Mini-batch stochastic approximation methods for nonconvex stochastic composite optimization. *Mathematical Programming*, 155(1):267–305, 2016. ISSN 1436-4646. doi: 10.1007/s10107-014-0846-1. URL <http://dx.doi.org/10.1007/s10107-014-0846-1>. arXiv:1308.6594.
- Elad Hazan and Kfir Levy. Bandit convex optimization: Towards tight bounds. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 784–792. Curran Associates, Inc., 2014. URL <http://papers.nips.cc/paper/5377-bandit-convex-optimization-towards-tight-bounds>.
- Xiaowei Hu, Prashanth L.A., Andrs Gyrgy, and Csaba Szepesvari. (bandit) convex optimization with biased noisy gradient oracles. In Arthur Gretton and Christian C. Robert, editors, *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, volume 51 of *Proceedings of Machine Learning Research*, pages 819–828, Cadiz, Spain, 09–11 May 2016. PMLR. URL <http://proceedings.mlr.press/v51/hu16b.html>.
- Kevin G Jamieson, Robert Nowak, and Ben Recht. Query complexity of derivative-free optimization. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 2672–2680. Curran Associates, Inc., 2012. URL <http://papers.nips.cc/paper/4509-query-complexity-of-derivative-free-optimization>.
- Guanghui Lan. An optimal method for stochastic composite optimization. *Mathematical Programming*, 133(1):365–397, Jun 2012. ISSN 1436-4646. URL <https://doi.org/10.1007/s10107-010-0434-y>. First appeared in June 2008.
- Tengyuan Liang, Hariharan Narayanan, and Alexander Rakhlin. On zeroth-order stochastic convex optimization via random walks. *arXiv:1402.2667*, 2014.
- A.S. Nemirovsky and D.B. Yudin. *Problem Complexity and Method Efficiency in Optimization*. J. Wiley & Sons, New York, 1983.
- Yurii Nesterov and Vladimir Spokoiny. Random gradient-free minimization of convex functions. *Found. Comput. Math.*, 17(2):527–566, April 2017. ISSN 1615-3375. doi: 10.1007/s10208-015-9296-2. URL <https://doi.org/10.1007/s10208-015-9296-2>. First appeared in 2011 as CORE discussion paper 2011/16.
- V. Yu. Protasov. Algorithms for approximate calculation of the minimum of a convex function from its values. *Mathematical Notes*, 59(1):69–74, Jan 1996. ISSN 1573-8876. doi: 10.1007/BF02312467. URL <https://doi.org/10.1007/BF02312467>.

- H. H. Rosenbrock. An automatic method for finding the greatest or least value of a function. *The Computer Journal*, 3(3):175–184, 1960. doi: 10.1093/comjnl/3.3.175. URL [+http://dx.doi.org/10.1093/comjnl/3.3.175](http://dx.doi.org/10.1093/comjnl/3.3.175).
- Ankan Saha and Ambuj Tewari. Improved regret guarantees for online smooth convex optimization with bandit feedback. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 636–642, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR. URL <http://proceedings.mlr.press/v15/sahaa11a.html>.
- Ohad Shamir. On the complexity of bandit and derivative-free stochastic convex optimization. In Shai Shalev-Shwartz and Ingo Steinwart, editors, *Proceedings of the 26th Annual Conference on Learning Theory*, volume 30 of *Proceedings of Machine Learning Research*, pages 3–24, Princeton, NJ, USA, 12–14 Jun 2013. PMLR. URL <http://proceedings.mlr.press/v30/Shamir13.html>.
- Ohad Shamir. An optimal algorithm for bandit and zero-order convex optimization with two-point feedback. *Journal of Machine Learning Research*, 18:52:1–52:11, 2017. URL <http://jmlr.org/papers/v18/papers/v18/16-632.html>. First appeared in arXiv:1507.08752.
- James C. Spall. *Introduction to Stochastic Search and Optimization*. John Wiley & Sons, Inc., New York, NY, USA, 1 edition, 2003. ISBN 0471330523.
- Sebastian Stich, Christian Myüller, and Bernd Gärtner. Optimization of convex functions with random pursuit. *arXiv:1111.0194*, 2011.

Appendix A. Proof of Lemma 1

Here we prove that, for $e \in RS_2(1)$

$$\mathbb{E}[\|e\|_q^2] \leq \min\{q-1, 16 \ln n - 8\} n^{\frac{2}{q}-1}, \quad (45)$$

$$\mathbb{E}[\langle s, e \rangle^2 \|e\|_q^2] \leq 6 \|s\|_2^2 \min\{q-1, 16 \ln n - 8\} n^{\frac{2}{q}-2}. \quad (46)$$

We start with proving the following inequality, which could be rough for big q :

$$\mathbb{E}[\|e\|_q^2] \leq (q-1) n^{\frac{2}{q}-1}, \quad 2 \leq q < \infty. \quad (47)$$

We have

$$\mathbb{E}[\|e\|_q^2] = \mathbb{E} \left[\left(\sum_{k=1}^n |e_k|^q \right)^{\frac{2}{q}} \right] \stackrel{\textcircled{1}}{\leq} \left(\mathbb{E} \left[\sum_{k=1}^n |e_k|^q \right] \right)^{\frac{2}{q}} \stackrel{\textcircled{2}}{=} (n \mathbb{E}[|e_2|^q])^{\frac{2}{q}}, \quad (48)$$

where $\textcircled{1}$ is due to probabilistic version of Jensen’s inequality (function $\varphi(x) = x^{\frac{2}{q}}$ is concave, because $q \geq 2$) and $\textcircled{2}$ is because mathematical expectation is linear and components of vector e are identically distributed.

Moreover, due to Poincare lemma, we have

$$e \stackrel{d}{=} \frac{\xi}{\sqrt{\xi_1^2 + \dots + \xi_n^2}}, \quad (49)$$

where ξ is Gaussian random vector which mathematical expectation is zero vector and covariance matrix is identical. Then

$$\begin{aligned} \mathbb{E}[|e_2|^q] &= \mathbb{E}\left[\frac{|\xi_2|^q}{(\xi_1^2 + \dots + \xi_n^2)^{\frac{q}{2}}}\right] \\ &= \int \dots \int_{\mathbb{R}^n} |x_2|^q \left(\sum_{k=1}^n x_k^2\right)^{-\frac{q}{2}} \cdot \frac{1}{(2\pi)^{\frac{n}{2}}} \cdot \exp\left(-\frac{1}{2} \sum_{k=1}^n x_k^2\right) dx_1 \dots dx_n. \end{aligned}$$

Consider spherical coordinates:

$$\begin{aligned} x_1 &= r \cos \varphi \sin \theta_1 \dots \sin \theta_{n-2}, \\ x_2 &= r \sin \varphi \sin \theta_1 \dots \sin \theta_{n-2}, \\ x_3 &= r \cos \theta_1 \sin \theta_2 \dots \sin \theta_{n-2}, \\ x_4 &= r \cos \theta_2 \sin \theta_3 \dots \sin \theta_{n-2}, \\ &\dots \\ x_n &= r \cos \theta_{n-2}, \\ r &> 0, \varphi \in [0, 2\pi), \theta_i \in [0, \pi], i = \overline{1, n-2}. \end{aligned}$$

The Jacobian of mapping is

$$\det \left(\frac{\partial(x_1, \dots, x_n)}{\partial(r, \varphi, \theta_1, \theta_2, \dots, \theta_{n-2})} \right) = r^{n-1} \sin \theta_1 (\sin \theta_2)^2 \dots (\sin \theta_{n-2})^{n-2}.$$

Then mathematical expectation $\mathbb{E}[|e_2|^q]$ could be rewritten in the following form:

$$\begin{aligned} &= \int \dots \int_{\substack{r>0, \varphi \in [0, 2\pi), \\ \theta_i \in [0, \pi], i=\overline{1, n-2}}} r^{n-1} |\sin \varphi|^q |\sin \theta_1|^{q+1} |\sin \theta_2|^{q+2} \dots |\sin \theta_{n-2}|^{q+n-2} \\ &\quad \cdot \frac{e^{-\frac{r^2}{2}}}{(2\pi)^{\frac{n}{2}}} dr \dots d\theta_{n-2} \\ &= \frac{1}{(2\pi)^{\frac{n}{2}}} I_r \cdot I_\varphi \cdot I_{\theta_1} \cdot I_{\theta_2} \cdot \dots \cdot I_{\theta_{n-2}}, \end{aligned}$$

where

$$\begin{aligned} I_r &= \int_0^{+\infty} r^{n-1} e^{-\frac{r^2}{2}} dr, \\ I_\varphi &= \int_0^{2\pi} |\sin \varphi|^q d\varphi = 2 \int_0^\pi |\sin \varphi|^q d\varphi, \\ I_{\theta_i} &= \int_0^\pi |\sin \theta_i|^{q+i} d\theta_i, i = \overline{1, n-2}. \end{aligned}$$

Now we are going to compute these integrals. Start with I_r :

$$I_r = \int_0^{+\infty} r^{n-1} e^{-\frac{r^2}{2}} dr = \int_0^{+\infty} (2t)^{\frac{n}{2}-1} e^{-t} dt = 2^{\frac{n}{2}-1} \Gamma\left(\frac{n}{2}\right).$$

To compute other integrals it is useful to consider the following integral ($\alpha > 0$):

$$\begin{aligned} \int_0^\pi |\sin \varphi|^\alpha d\varphi &= 2 \int_0^{\frac{\pi}{2}} |\sin \varphi|^\alpha d\varphi = 2 \int_0^{\frac{\pi}{2}} (\sin^2 \varphi)^{\frac{\alpha}{2}} d\varphi = \int_0^1 t^{\frac{\alpha-1}{2}} (1-t)^{-\frac{1}{2}} dt = B\left(\frac{\alpha+1}{2}, \frac{1}{2}\right) = \frac{\Gamma(\frac{\alpha+1}{2})\Gamma(\frac{1}{2})}{\Gamma(\frac{\alpha+2}{2})} = \sqrt{\pi} \frac{\Gamma(\frac{\alpha+1}{2})}{\Gamma(\frac{\alpha+2}{2})}. \end{aligned}$$

From this we obtain

$$\begin{aligned} \mathbb{E}[|e_2|^q] &= \frac{1}{(2\pi)^{\frac{n}{2}}} I_r \cdot I_\varphi \cdot I_{\theta_1} \cdot I_{\theta_2} \cdot \dots \cdot I_{\theta_{n-2}} \\ &= \frac{1}{(2\pi)^{\frac{n}{2}}} \cdot 2^{\frac{n}{2}-1} \Gamma\left(\frac{n}{2}\right) \cdot 2\sqrt{\pi} \frac{\Gamma(\frac{q+1}{2})}{\Gamma(\frac{q+2}{2})} \cdot \sqrt{\pi} \frac{\Gamma(\frac{q+2}{2})}{\Gamma(\frac{q+3}{2})} \cdot \sqrt{\pi} \frac{\Gamma(\frac{q+3}{2})}{\Gamma(\frac{q+4}{2})} \cdot \dots \cdot \sqrt{\pi} \frac{\Gamma(\frac{q+n-1}{2})}{\Gamma(\frac{q+n}{2})} \\ &= \frac{1}{\sqrt{\pi}} \cdot \frac{\Gamma(\frac{n}{2})\Gamma(\frac{q+1}{2})}{\Gamma(\frac{q+n}{2})}. \end{aligned} \quad (50)$$

Now, we want to show that $\forall q \geq 2$

$$\frac{1}{\sqrt{\pi}} \cdot \frac{\Gamma(\frac{n}{2})\Gamma(\frac{q+1}{2})}{\Gamma(\frac{q+n}{2})} \leq \left(\frac{q-1}{n}\right)^{\frac{q}{2}}. \quad (51)$$

At the beginning show that (51) holds for $q = 2$ (and arbitrary n):

$$\frac{1}{\sqrt{\pi}} \cdot \frac{\Gamma(\frac{n}{2})\Gamma(\frac{2+1}{2})}{\Gamma(\frac{2+n}{2})} - \frac{1}{n} = \frac{1}{\sqrt{\pi}} \cdot \frac{\Gamma(\frac{n}{2}) \cdot \frac{1}{2}\Gamma(\frac{1}{2})}{\frac{n}{2}\Gamma(\frac{n}{2})} - \frac{1}{n} = \frac{1}{n} - \frac{1}{n} = 0 \leq 0.$$

Consider the function

$$f_n(q) = \frac{1}{\sqrt{\pi}} \cdot \frac{\Gamma(\frac{n}{2})\Gamma(\frac{q+1}{2})}{\Gamma(\frac{q+n}{2})} - \left(\frac{q-1}{n}\right)^{\frac{q}{2}}$$

where $q \geq 2$. Also consider $\psi(x) = \frac{d(\ln(\Gamma(x)))}{dx}$ with $x > 0$ which is called (*digamma function*). For gamma function it holds

$$\Gamma(x+1) = x\Gamma(x), \quad x > 0.$$

Taking natural logarithm from it and taking derivative w.r.t. x :

$$\begin{aligned} \ln \Gamma(x+1) &= \ln \Gamma(x) + \ln x, \\ \frac{d(\ln(\Gamma(x+1)))}{dx} &= \frac{d(\ln(\Gamma(x)))}{dx} + \frac{1}{x}, \end{aligned}$$

which could be written in digamma-function-notation:

$$\psi(x+1) = \psi(x) + \frac{1}{x}. \quad (52)$$

One can show that digamma function is monotonically increases when $x > 0$. To prove this fact we are going to show that

$$(\Gamma'(x))^2 < \Gamma(x)\Gamma''(x). \quad (53)$$

That is,

$$\begin{aligned} (\Gamma'(x))^2 &= \left(\int_0^{+\infty} e^{-t} \ln t \cdot t^{x-1} dt \right)^2 \\ &\stackrel{\textcircled{1}}{<} \int_0^{+\infty} \left(e^{-\frac{t}{2}} t^{\frac{x-1}{2}} \right)^2 dt \cdot \int_0^{+\infty} \left(e^{-\frac{t}{2}} t^{\frac{x-1}{2}} \ln t \right)^2 dt = \underbrace{\int_0^{+\infty} e^{-t} t^{x-1} dt}_{\Gamma(x)} \cdot \underbrace{\int_0^{+\infty} e^t t^{x-1} \ln^2 t dt}_{\Gamma''(x)}, \end{aligned}$$

where $\textcircled{1}$ follows from Cauchy-Schwartz inequality (the equality cannot occur because functions $e^{-\frac{t}{2}} t^{\frac{x-1}{2}}$ and $e^{-\frac{t}{2}} t^{\frac{x-1}{2}} \ln t$ are linearly independent). From (53) follows that

$$\frac{d^2(\ln \Gamma(x))}{dx^2} = \left(\frac{\Gamma'(x)}{\Gamma(x)} \right)' = \frac{\Gamma''(x)}{\Gamma(x)} - \frac{(\Gamma'(x))^2}{(\Gamma(x))^2} \stackrel{(53)}{>} 0,$$

which shows that digamma function increases.

Now we show that $f_n(q)$ decreases on the interval $[2, +\infty)$. To obtain it is sufficient to consider $\ln(f(q))$:

$$\begin{aligned} &\ln(f_n(q)) \\ &= \ln \left(\frac{\Gamma(\frac{n}{2})}{\sqrt{\pi}} \right) + \ln \left(\Gamma \left(\frac{q+1}{2} \right) \right) - \ln \left(\Gamma \left(\frac{q+n}{2} \right) \right) - \frac{q}{2} (\ln(q-1) - \ln n), \\ \frac{d(\ln(f_n(q)))}{dq} &= \frac{1}{2} \psi \left(\frac{q+1}{2} \right) - \frac{1}{2} \psi \left(\frac{q+n}{2} \right) - \frac{1}{2} \ln(q-1) - \frac{q}{2(q-1)} + \frac{1}{2} \ln n. \end{aligned}$$

We are going to show that $\frac{d(\ln(f_n(q)))}{dq} < 0$ for $q \geq 2$. Let $k = \lfloor \frac{n}{2} \rfloor$ (the closest integer which is no greater than $\frac{n}{2}$). Then $\psi \left(\frac{q+n}{2} \right) > \psi \left(k-1 + \frac{q+1}{2} \right)$ and $\ln n \leq \ln(2k+1)$, whence

$$\begin{aligned} &\frac{d(\ln(f_n(q)))}{dq} \\ &< \frac{1}{2} \left(\psi \left(\frac{q+1}{2} \right) - \psi \left(k-1 + \frac{q+1}{2} \right) \right) - \frac{1}{2} \ln(q-1) - \frac{q}{2(q-1)} + \frac{1}{2} \ln(2k+1) \\ &\stackrel{(52)}{=} \frac{1}{2} \left(\psi \left(\frac{q+1}{2} \right) - \sum_{i=1}^{k-1} \frac{1}{\frac{q+1}{2} + k - i - 1} - \psi \left(\frac{q+1}{2} \right) \right) - \frac{q}{2(q-1)} + \frac{1}{2} \ln \left(\frac{2k+1}{q-1} \right) \\ &\stackrel{\textcircled{1}}{\leq} -\frac{1}{2} \sum_{i=1}^{k-1} \frac{2}{q-1+2k-2i} - \frac{1}{q-1} + \frac{1}{2} \ln \left(\frac{2k+1}{q-1} \right) \\ &= -\frac{1}{2} \left(\frac{2}{q-1} + \frac{2}{q+1} + \frac{2}{q+3} + \dots + \frac{2}{q+2k-3} \right) + \frac{1}{2} \ln \left(\frac{2k+1}{q-1} \right) \\ &\stackrel{\textcircled{2}}{<} -\frac{1}{2} \ln \left(\frac{q+2k-1}{q-1} \right) + \frac{1}{2} \ln \left(\frac{2k+1}{q-1} \right) \stackrel{\textcircled{3}}{\leq} -\frac{1}{2} \ln \left(\frac{2k+1}{q-1} \right) + \frac{1}{2} \ln \left(\frac{2k+1}{q-1} \right) = 0, \end{aligned}$$

where $\textcircled{1}$ and $\textcircled{3}$ is because $q \geq 2$, $\textcircled{2}$ is due to estimation of integral of $\frac{1}{x}$ by integral of $g(x) = \frac{1}{q-1+2i}$, $x \in [q-1+2i, q-1+2i+2]$, $i = \overline{0, 2k-1}$ which is no less than $f(x)$:

$$\frac{2}{q-1} + \frac{2}{q+1} + \frac{2}{q+3} + \dots + \frac{2}{q+2k-3} > \int_{q-1}^{q+2k-1} \frac{1}{x} dx = \ln \left(\frac{q+2k-1}{q-1} \right).$$

So, we shown that $\frac{d(\ln(f_n(q)))}{dq} < 0$ for $q \geq 2$ arbitrary natural number n . Therefore for any fixed number n the function $f_n(q)$ decreases as q increase, which means that $f_n(q) \leq f_n(2) = 0$, i.e., (51) holds. From this and (48),(50) we obtain that $\forall q \geq 2$

$$\mathbb{E}[\|e\|_q^2] \stackrel{(48)}{\leq} (n\mathbb{E}[|e_2|^q])^{\frac{2}{q}} \stackrel{(50),(51)}{\leq} (q-1)n^{\frac{2}{q}-1}. \quad (54)$$

However, inequality (54) is useless when q is big (with respect to n). Consider left hand side of (54) as function of q and find its minimum for $q \geq 2$. Consider $h_n(q) = \ln(q-1) + \left(\frac{2}{q}-1\right) \ln n$ (it is logarithm of the right hand side of (54)). Derivative of $h(q)$ is

$$\begin{aligned} \frac{dh(q)}{dq} &= \frac{1}{q-1} - \frac{2\ln n}{q^2}, \\ \frac{1}{q-1} - \frac{2\ln n}{q^2} &= 0, \\ q^2 - 2q \ln n + 2 \ln n &= 0. \end{aligned}$$

If $n \geq 8$, then the point where the function obtains its minimum on the set $[2, +\infty)$ is $q_0 = \ln n \left(1 + \sqrt{1 - \frac{2}{\ln n}}\right)$ (for the case $n \leq 7$ it turns out that $q_0 = 2$; further without loss of generality we assume $n \geq 8$). Therefore for all $q > q_0$ it is more useful to use the following estimation:

$$\begin{aligned} \mathbb{E}[\|e\|_q^2] &\stackrel{\textcircled{1}}{<} \mathbb{E}[\|e\|_{q_0}^2] \stackrel{(54)}{\leq} (q_0-1)n^{\frac{2}{q_0}-1} \stackrel{\textcircled{2}}{\leq} (2\ln n-1)n^{\frac{2}{\ln n}-1} \\ &= (2\ln n-1)e^2 \frac{1}{n} \leq (16\ln n-8) \frac{1}{n} \leq (16\ln n-8)n^{\frac{2}{q}-1}, \end{aligned} \quad (55)$$

where $\textcircled{1}$ is due to $\|e\|_q < \|e\|_{q_0}$ for $q > q_0$, $\textcircled{2}$ follows from $q_0 \leq 2\ln n$, $q_0 \geq \ln n$. Putting estimations (54) and (55) together we obtain (45).

Now we are going to prove (46). Firstly, we want to estimate $\sqrt{\mathbb{E}[\|e\|_q^4]}$. Due to probabilistic Jensen's inequality ($q \geq 2$)

$$\begin{aligned} \mathbb{E}[\|e\|_q^4] &= \mathbb{E} \left[\left(\left(\sum_{k=1}^n |e_k|^q \right)^2 \right)^{\frac{2}{q}} \right] \leq \left(\mathbb{E} \left[\left(\sum_{k=1}^n |e_k|^q \right)^2 \right] \right)^{\frac{2}{q}} \\ &\stackrel{\textcircled{1}}{\leq} \left(\mathbb{E} \left[\left(n \sum_{k=1}^n |e_k|^{2q} \right) \right] \right)^{\frac{2}{q}} \stackrel{\textcircled{2}}{=} (n^2 \mathbb{E}[|e_2|^{2q}])^{\frac{2}{q}} \\ &\stackrel{(50),(51)}{\leq} n^{\frac{4}{q}} \left(\left(\frac{2q-1}{n} \right)^{\frac{2q}{2}} \right)^{\frac{2}{q}} = (2q-1)^2 n^{\frac{4}{q}-2}, \end{aligned}$$

where $\textcircled{1}$ is because $\left(\sum_{k=1}^n x_k \right)^2 \leq n \sum_{k=1}^n x_k^2$ for $x_1, x_2, \dots, x_n \in \mathbb{R}$ and $\textcircled{2}$ follows from that mathematical expectation is linear and components of the random vector e are identically distributed. From this we obtain

$$\sqrt{\mathbb{E}[\|e\|_q^4]} \leq (2q-1)n^{\frac{2}{q}-1}. \quad (56)$$

Consider the right hand side of the inequality (56) as a function of q and find its minimum for $q \geq 2$. Consider $h_n(q) = \ln(2q-1) + \left(\frac{2}{q}-1\right) \ln n$ (logarithm of the right hand side (56)). Derivative of $h(q)$ is

$$\begin{aligned} \frac{dh(q)}{dq} &= \frac{2}{2q-1} - \frac{2\ln n}{q^2}, \\ \frac{2}{2q-1} - \frac{2\ln n}{q^2} &= 0, \\ q^2 - 2q \ln n + \ln n &= 0. \end{aligned}$$

If $n \geq 3$, the the point where the function obtains its minimum on the set $[2, +\infty)$ is $q_0 = \ln n \left(1 + \sqrt{1 - \frac{1}{\ln n}}\right)$ (for the case $n \leq 2$ it turns out that $q_0 = 2$; further without loss of generality we assume that $n \geq 3$). Therefore for all $q > q_0$:

$$\begin{aligned} \sqrt{\mathbb{E}[||e||_q^4]} &\stackrel{\textcircled{1}}{\leq} \sqrt{\mathbb{E}[||e||_{q_0}^4]} \stackrel{(56)}{\leq} (2q_0 - 1)n^{\frac{2}{q_0}-1} \stackrel{\textcircled{2}}{\leq} (4\ln n - 1)n^{\frac{2}{\ln n}-1} \\ &= (4\ln n - 1)e^2 \frac{1}{n} \leq (32\ln n - 8) \frac{1}{n} \leq (32\ln n - 8)n^{\frac{2}{q}-1}, \end{aligned} \quad (57)$$

where $\textcircled{1}$ is due to $||e||_q < ||e||_{q_0}$ for $q > q_0$, $\textcircled{2}$ follows from $q_0 \leq 2\ln n$, $q_0 \geq \ln n$. Putting estimations (56) and (57) together we get inequality

$$\sqrt{\mathbb{E}[||e||_q^4]} \leq \min\{2q - 1, 32\ln n - 8\}n^{\frac{2}{q}-1}. \quad (58)$$

Now we are going to find $\mathbb{E}[\langle s, e \rangle^4]$, where $s \in \mathbb{R}^n$ is some vector. Let $S_n(r)$ be a surface area of n -dimensional Euclidean sphere with radius r and $d\sigma(e)$ be unnormalized uniform measure on n -dimensional Euclidean sphere. From this it follows that $S_n(r) = S_n(1)r^{n-1}$, $\frac{S_{n-1}(1)}{S_n(1)} = \frac{n-1}{n\sqrt{\pi}} \frac{\Gamma(\frac{n+2}{2})}{\Gamma(\frac{n+1}{2})}$. Besides, let φ be the angle between s and e . Then

$$\begin{aligned} \mathbb{E}[\langle s, e \rangle^4] &= \frac{1}{S_n(1)} \int_S \langle s, e \rangle^4 d\sigma(\varphi) = \frac{1}{S_n(1)} \int_0^\pi ||s||_2^4 \cos^3 \varphi S_{n-1}(\sin \varphi) d\varphi \\ &= ||s||_2^4 \frac{S_{n-1}(1)}{S_n(1)} \int_0^\pi \cos^4 \varphi \sin^{n-2} \varphi d\varphi = ||s||_2^4 \cdot \frac{n-1}{n\sqrt{\pi}} \frac{\Gamma(\frac{n+2}{2})}{\Gamma(\frac{n+1}{2})} \int_0^\pi \cos^4 \varphi \sin^{n-2} \varphi d\varphi. \end{aligned} \quad (59)$$

Compute the integral:

$$\begin{aligned} \int_0^\pi \cos^4 \varphi \sin^{n-2} \varphi d\varphi &= 2 \int_0^{\frac{\pi}{2}} \cos^4 \varphi \sin^{n-2} \varphi d\varphi = \int_0^{\frac{\pi}{2}} t^{\frac{n-3}{2}} (1-t)^{\frac{3}{2}} dt = B\left(\frac{n-1}{2}, \frac{5}{2}\right) = \frac{\Gamma(\frac{5}{2})\Gamma(\frac{n-1}{2})}{\Gamma(\frac{n+4}{2})} = \frac{\frac{3}{2} \cdot \frac{1}{2} \Gamma(\frac{1}{2}) \Gamma(\frac{n-1}{2})}{\frac{n+2}{2} \Gamma(\frac{n+2}{2})} = \frac{3}{n+2} \cdot \frac{\sqrt{\pi} \Gamma(\frac{n-1}{2})}{2\Gamma(\frac{n+2}{2})}. \end{aligned}$$

From this and (59) we obtain

$$\begin{aligned} \mathbb{E}[\langle s, e \rangle^4] &= ||s||_2^4 \cdot \frac{n-1}{n\sqrt{\pi}} \frac{\Gamma(\frac{n+2}{2})}{\Gamma(\frac{n+1}{2})} \cdot \frac{3}{n+2} \cdot \frac{\sqrt{\pi} \Gamma(\frac{n-1}{2})}{2\Gamma(\frac{n+2}{2})} \\ &= ||s||_2^4 \cdot \frac{3(n-1)}{2n(n+2)} \cdot \frac{\Gamma(\frac{n-1}{2})}{\frac{n-1}{2} \Gamma(\frac{n-1}{2})} = \frac{3||s||_2^4}{n(n+2)} \stackrel{\textcircled{1}}{\leq} \frac{3||s||_2^4}{n^2}. \end{aligned} \quad (60)$$

To prove (46), it remains to use (58), (60) and Cauchy-Schwartz inequality $(\mathbb{E}[XY])^2 \leq \mathbb{E}[X^2] \cdot \mathbb{E}[Y^2]$:

$$\mathbb{E}[\langle s, e \rangle^2 ||e||_q^2] \stackrel{\textcircled{1}}{\leq} \sqrt{\mathbb{E}[\langle s, e \rangle^4] \cdot \mathbb{E}[||e||_q^4]} \leq \sqrt{3} ||s||_2^2 \min\{2q - 1, 32\ln n - 8\}n^{\frac{2}{q}-2}.$$

Appendix B. Technical Results

Lemma 8 *Let $a_0, \dots, a_{N-1}, b, R_1, \dots, R_{N-1}$ be non-negative numbers such that*

$$R_l \leq \sqrt{2} \cdot \sqrt{\left(\sum_{k=0}^{l-1} a_k + b \sum_{k=1}^{l-1} \alpha_{k+1} R_k \right)} \quad l = 1, \dots, N, \quad (61)$$

where $\alpha_{k+1} = \frac{k+2}{96n^2\rho_n L_2}$ for all $k \in \mathbb{N}$. Then for $l = 1, \dots, N$

$$\sum_{k=0}^{l-1} a_k + b \sum_{k=1}^{l-1} \alpha_{k+1} R_k \leq \left(\sqrt{\sum_{k=0}^{l-1} a_k} + \sqrt{2b} \cdot \frac{l^2}{96n^2\rho_n L_2} \right)^2. \quad (62)$$

Proof For $l = 1$ it is trivial inequality. Assume that (62) holds for some $l < N$ and prove it for $l + 1$. From the induction assumption and (61) we obtain

$$R_l \leq \sqrt{2} \left(\sqrt{\sum_{k=0}^{l-1} a_k} + \sqrt{2b} \cdot \frac{l^2}{96n^2\rho_n L_2} \right), \quad (63)$$

whence

$$\begin{aligned} \sum_{k=0}^l a_k + b \sum_{k=1}^l \alpha_{k+1} R_k &= \sum_{k=0}^{l-1} a_k + b \sum_{k=1}^{l-1} \alpha_{k+1} R_k + a_l + b\alpha_{l+1} R_l \\ &\stackrel{\textcircled{1}}{\leq} \left(\sqrt{\sum_{k=0}^{l-1} a_k} + \sqrt{2b} \cdot \frac{l^2}{96n^2\rho_n L_2} \right)^2 + a_l + \sqrt{2b}\alpha_{l+1} \left(\sqrt{\sum_{k=0}^{l-1} a_k} + \sqrt{2b} \cdot \frac{l^2}{96n^2\rho_n L_2} \right) \\ &= \sum_{k=0}^l a_k + 2\sqrt{\sum_{k=0}^{l-1} a_k} \cdot \sqrt{2b} \frac{l^2}{96n^2\rho_n L_2} + 2b^2 \frac{l^4}{(96n^2\rho_n L_2)^2} + \sqrt{2b}\alpha_{l+1} \left(\sqrt{\sum_{k=0}^{l-1} a_k} + \sqrt{2b} \cdot \frac{l^2}{96n^2\rho_n L_2} \right) \\ &= \sum_{k=0}^l a_k + 2\sqrt{\sum_{k=0}^{l-1} a_k} \cdot \sqrt{2b} \left(\frac{l^2}{96n^2\rho_n L_2} + \frac{\alpha_{l+1}}{2} \right) + 2b^2 \left(\frac{l^4}{(96n^2\rho_n L_2)^2} + \alpha_{l+1} \cdot \frac{l^2}{96n^2\rho_n L_2} \right) \\ &\stackrel{\textcircled{2}}{\leq} \sum_{k=0}^l a_k + 2\sqrt{\sum_{k=0}^l a_k} \cdot \sqrt{2b} \frac{(l+1)^2}{96n^2\rho_n L_2} + 2b^2 \frac{(l+1)^4}{(96n^2\rho_n L_2)^2} = \left(\sqrt{\sum_{k=0}^l a_k} + \sqrt{2b} \cdot \frac{(l+1)^2}{96n^2\rho_n L_2} \right)^2, \end{aligned}$$

where $\textcircled{1}$ follows from the induction assumption and (63), $\textcircled{2}$ is because $\sum_{k=0}^{l-1} a_k \leq \sum_{k=0}^l a_k$ and

$$\begin{aligned} \frac{l^2}{96n^2\rho_n L_2} + \frac{\alpha_{l+1}}{2} &= \frac{2l^2 + l + 2}{192n^2\rho_n L_2} \leq \frac{(l+1)^2}{96n^2\rho_n L_2}, \\ \frac{l^4}{(96n^2\rho_n L_2)^2} + \alpha_{l+1} \cdot \frac{l^2}{96n^2\rho_n L_2} &\leq \frac{l^4 + (l+2)l^2}{(96n^2\rho_n L_2)^2} \leq \frac{(l+1)^4}{(96n^2\rho_n L_2)^2}. \end{aligned}$$

■

Lemma 9 Let $a_0, \dots, a_{N-1}, b, R_1, \dots, R_{N-1}$ be non-negative numbers such that

$$R_l \leq \sqrt{2} \cdot \sqrt{\left(\sum_{k=0}^{l-1} a_k + b\alpha \sum_{k=1}^{l-1} R_k \right)} \quad l = 1, \dots, N. \quad (64)$$

Then for $l = 1, \dots, N$

$$\sum_{k=0}^{l-1} a_k + b\alpha \sum_{k=1}^{l-1} R_k \leq \left(\sqrt{\sum_{k=0}^{l-1} a_k} + \sqrt{2b\alpha l} \right)^2. \quad (65)$$

Proof For $l = 1$ it is trivial inequality. Assume that (65) holds for some $l < N$ and prove it for $l + 1$. From the induction assumption and (64) we obtain

$$R_l \leq \sqrt{2} \left(\sqrt{\sum_{k=0}^{l-1} a_k} + \sqrt{2b\alpha l} \right), \quad (66)$$

whence

$$\begin{aligned} & \sum_{k=0}^l a_k + b\alpha \sum_{k=1}^l R_k = \sum_{k=0}^{l-1} a_k + b\alpha \sum_{k=1}^{l-1} R_k + a_l + b\alpha R_l \\ & \stackrel{\textcircled{1}}{\leq} \left(\sqrt{\sum_{k=0}^{l-1} a_k} + \sqrt{2b\alpha l} \right)^2 + a_l + \sqrt{2b\alpha} \left(\sqrt{\sum_{k=0}^{l-1} a_k} + \sqrt{2b\alpha l} \right) \\ & = \sum_{k=0}^l a_k + 2\sqrt{\sum_{k=0}^{l-1} a_k} \cdot \sqrt{2b\alpha l} + 2b^2\alpha^2 l^2 + \sqrt{2b\alpha} \left(\sqrt{\sum_{k=0}^{l-1} a_k} + \sqrt{2b\alpha l} \right) \\ & = \sum_{k=0}^l a_k + 2\sqrt{\sum_{k=0}^{l-1} a_k} \cdot \sqrt{2b\alpha} \left(l + \frac{1}{2} \right) + 2b^2\alpha^2 (l^2 + l) \\ & \stackrel{\textcircled{2}}{\leq} \sum_{k=0}^l a_k + 2\sqrt{\sum_{k=0}^l a_k} \cdot \sqrt{2b\alpha} (l + 1) + 2b^2\alpha^2 (l + 1)^2 = \left(\sqrt{\sum_{k=0}^l a_k} + \sqrt{2b\alpha} (l + 1) \right)^2, \end{aligned}$$

where $\textcircled{1}$ follows from the induction assumption and (66), $\textcircled{2}$ is because $\sum_{k=0}^{l-1} a_k \leq \sum_{k=0}^l a_k$. ■