

# Randomized first order algorithms with applications to $\ell_1$ -minimization

Anatoli Juditsky · Fatma Kılınç Karzan ·  
Arkadi Nemirovski

Received: 8 February 2011 / Accepted: 26 June 2012 / Published online: 25 July 2012  
© Springer and Mathematical Optimization Society 2012

**Abstract** In this paper we propose randomized first-order algorithms for solving bilinear saddle points problems. Our developments are motivated by the need for sub-linear time algorithms to solve large-scale *parametric* bilinear saddle point problems where cheap online assessment of the solution quality is crucial. We present the theoretical efficiency estimates of our algorithms and discuss a number of applications, primarily to the problem of  $\ell_1$  minimization arising in sparsity-oriented signal processing. We demonstrate, both theoretically and by numerical examples, that when seeking for medium-accuracy solutions of large-scale  $\ell_1$  minimization problems, our randomized algorithms outperform significantly (and progressively as the sizes of the problem grow) the state-of-the art deterministic methods.

**Mathematics Subject Classification** 90C25 · 90C47 · 90C06 · 65K15

## 1 Introduction

This paper is motivated by the desire to develop efficient *randomized* first-order methods for solving well-structured large-scale convex optimization problems.

---

Research was partly supported by the ONR grant N000140811104 and the NSF grant DMS-0914785 (second and third authors) and the BSF grant 2008302 (third author).

---

A. Juditsky  
LJK, Université J. Fourier, B.P. 53, 38041 Grenoble Cedex 9, France  
e-mail: Anatoli.Juditsky@imag.fr

F. Kılınç Karzan (✉)  
Carnegie Mellon University, Pittsburgh, PA 15213, USA  
e-mail: fkilinc@andrew.cmu.edu

A. Nemirovski  
Georgia Institute of Technology, Atlanta, GA 30332, USA  
e-mail: nemirovs@isye.gatech.edu

Our approach is based on saddle point (s.p.) reformulation of well-structured convex minimization problems and is applicable when the resulting s.p. problems are bilinear; in this respect, it goes back to the breakthrough paper of Nesterov [15]. The deterministic s.p. prototypes of the randomized algorithms we develop here were proposed in [12, 13], and the prototypes of our randomization scheme were proposed in [14, Section 3.3] and [10]. In this paper, we demonstrate that in the case of a bilinear s.p. problem, a better randomization is possible,<sup>1</sup> specifically, one allowing to assess in a computationally cheap fashion the quality of the resulting approximate solutions. This assessment is instrumental when solving parametric bilinear s.p. problems covering numerous applications.

As an application area, our primary (but not the only) target is the  $\ell_1$ -minimization problem

$$\text{Opt}_p = \min_u \{ \|u\|_1 : \|Au - b\|_p \leq \delta \} \quad [A = [A_1, \dots, A_n] \in \mathbf{R}^{m \times n}, m, n > 2], \quad (1)$$

where  $p = \infty$  (“uniform fit”) or  $p = 2$  (“ $\ell_2$ -fit”). We are interested in the large-scale case, where the sizes  $m, n$  of (possibly dense) matrix  $A$  are in the range of thousands/tens of thousands. Efficient solutions to the problems of this type are of paramount importance for sparsity-oriented Signal Processing, in particular, in compressed sensing (see [2, 3, 5] and references therein). To give a flavor of our results, here is what our approach yields for (1):

**Proposition 1** *Assume that (1) is feasible,  $\delta$  is small enough, namely,  $2m^{\frac{1}{p}}\delta \leq \|b\|_p$ . Given  $\epsilon \in (0, \frac{1}{2}\text{Opt}_p\|A\|_{1,p})$ ,<sup>2</sup> let our goal be to find an  $\epsilon$ -solution to (1), that is, a point  $x_\epsilon$  satisfying*

$$\|x_\epsilon\|_1 \leq \text{Opt}_p \ \& \ \|Ax_\epsilon - b\|_p \leq \delta + \epsilon.$$

*Then, for every tolerance  $\chi \in (0, 1/2]$ , the outlined goal can be achieved with probability  $\geq 1 - \chi$*

(i) *in the case of  $p = \infty$  (uniform fit)—in at most*

$$O(1) \left[ \frac{\sqrt{\ln(m)\ln(n)}\|A\|_{1,\infty}\text{Opt}_\infty}{\epsilon} \ln \left( \frac{\sqrt{\ln(m)\ln(n)}\|A\|_{1,\infty}\text{Opt}_\infty}{\chi\epsilon} \right) \right]^2$$

*steps of a randomized algorithm, with computational effort per step reduced to extracting from  $A$  two columns and two rows, given their indexes, plus “computational overhead” of  $O(1)(m + n)$  operations.*

<sup>1</sup> In the hindsight, a particular case of this new randomization can be recognized in the sublinear time randomized algorithm for matrix games due to Grigoriadis and Khachiyan [7].

<sup>2</sup> Here and below  $\|A\|_{1,p} = \max_j \|A_j\|_p$  is the norm of the mapping  $x \mapsto Ax$  induced by the norms  $\|\cdot\|_1$  and  $\|\cdot\|_p$  in the argument and the image spaces, respectively.

(ii) in the case of  $p = 2$  ( $\ell_2$  fit)—in at most

$$O(1) \left[ \frac{\ln(mn)\Gamma(A)\|A\|_{1,2}\text{Opt}_2}{\epsilon} \ln \left( \frac{\ln(mn)\Gamma(A)\|A\|_{1,2}\text{Opt}_2}{\chi\epsilon} \right) \right]^2,$$

$$\Gamma(A) = \frac{\sqrt{m}\|A\|_{1,\infty}}{\|A\|_{1,2}},$$

steps of a randomized algorithm with the same computational effort per step as in (i).

Furthermore, there exists a randomized preprocessing of the data  $[A, b]$  of the problem (1) of computational cost not exceeding  $O(1)mn \ln(m)$ , which ensures with probability  $\geq 1 - \chi$  that  $\Gamma(A) \leq O(1)\sqrt{\ln(mn/\chi)}$ .

Note that the best known so far complexity of finding  $\epsilon$ -solution to a large-scale problem (1) by a deterministic algorithm is at least  $O(1)\frac{\sqrt{\ln(m)\ln(n)}\|A\|_{1,\infty}\text{Opt}_\infty}{\epsilon}$  ( $p = \infty$ ) or  $O(1)\frac{\sqrt{\ln(n)}\|A\|_{1,2}\text{Opt}_2}{\epsilon}$  ( $p = 2$ ) steps<sup>3</sup> with complexity of a step dominated by the necessity to perform  $O(1)$  multiplications  $x \mapsto Ax, y \mapsto A^T y$ . When  $A$  is dense, full matrix vector product requires  $O(mn)$  operations, and hence the total operations count is, up to log-factors, of order of  $\frac{mn}{\nu}$ , where  $\nu = \frac{\epsilon}{\|A\|_{1,p}\text{Opt}_p}$  can be naturally interpreted as relative accuracy. For the randomized algorithms underlying Proposition 1, this count, again up to log-factors, is of order of  $\frac{m+n}{\nu^2}$  (uniform fit) and  $\frac{m+n}{\nu^2} + mn$  ( $\ell_2$  fit). We see that when the relative accuracy  $\nu$  is such that  $1 \gg \nu \gg m^{-1} + n^{-1}$ , the randomized algorithms outperform the deterministic ones, and the positive effect of randomization becomes more significant as the problem size grows, i.e., “ $\gg$ ” in the above becomes “sharper”. Numerical results presented in Sect. 5 demonstrate that this acceleration can be of real practical interest.

The main body of this paper is organized as follows. In Sect. 2, we present a saddle-point-based framework for our developments together with a sample of interesting optimization problems fitting this framework. This sample includes, along with  $\ell_1$  minimization, the (semidefinite relaxation of the) problem of low-dimensional approximation of a collection of points in  $\mathbf{R}^d$ . Randomized algorithms for the problems fitting to our framework are developed and analyzed in Sects. 3 and 4. Section 5 presents encouraging preliminary results on numerical comparison of our randomized algorithms and their state-of-the-art deterministic counterparts as applied to large-scale  $\ell_1$  minimization problems.

## 2 Problems and goals

We start with specifying and motivating two problems to be discussed in the paper and our goals.

<sup>3</sup> The bounds are attainable, provided  $\text{Opt}_p$  is known in advance.

### 2.1 A bilinear saddle point problem

The first problem we are interested in is a Bilinear Saddle Point (BSP) problem

$$\begin{aligned}
 SV &= \min_{z_1 \in Z_1} \max_{z_2 \in Z_2} \phi(z_1, z_2), \\
 \phi(z_1, z_2) &= v + \langle a_1, z_1 \rangle + \langle a_2, z_2 \rangle + \langle z_2, Bz_1 \rangle : Z := Z_1 \times Z_2 \rightarrow \mathbf{R},
 \end{aligned} \tag{S}$$

where  $Z_i$  are nonempty convex compact sets in Euclidean spaces  $E_i, i = 1, 2$ . Recall that (S) gives rise to two convex optimization programs that are dual to each other:

$$\begin{aligned}
 \text{Opt}(P) &= \min_{z_1 \in Z_1} \left[ \bar{\phi}(z_1) := \max_{z_2 \in Z_2} \phi(z_1, z_2) \right] \quad (P) \\
 \text{Opt}(D) &= \max_{z_2 \in Z_2} \left[ \underline{\phi}(z_2) := \min_{z_1 \in Z_1} \phi(z_1, z_2) \right] \quad (D)
 \end{aligned} \tag{2}$$

with  $\text{Opt}(P) = \text{Opt}(D) = SV$ , and to the variational inequality (v.i.):

$$\text{find } z_* \in Z := Z_1 \times Z_2 \text{ such that } \langle F(z), z - z_* \rangle \geq 0 \text{ for all } z \in Z, \tag{3}$$

where  $F : Z \mapsto E_1 \times E_2$  is the affine monotone operator given by

$$\begin{aligned}
 F(z_1, z_2) &= \left[ F_1(z_2) = \frac{\partial \phi(z_1, z_2)}{\partial z_1}; F_2(z_1) = -\frac{\partial \phi(z_1, z_2)}{\partial z_2} \right] = a + \mathcal{A}[z_1; z_2], \\
 a &= [a_1; -a_2], \quad \mathcal{A} = \left[ \begin{array}{c|c} & B^* \\ \hline -B & \end{array} \right]
 \end{aligned} \tag{4}$$

(here  $B^*$  stands for the conjugate of  $B$ ). Note that  $\mathcal{A}$  is skew-symmetric:

$$\langle z, \mathcal{A}z \rangle = 0 \quad \forall z \in E := E_1 \times E_2. \tag{5}$$

It is well known that the solutions to (S)—the saddle points of  $\phi$  on  $Z_1 \times Z_2$ —are exactly the pairs  $z = [z_1; z_2]$  comprised of optimal solutions to problems (P) and (D) in (2). They are also exactly the solutions to the v.i. (3). We quantify the accuracy of candidate solutions  $z = [z_1; z_2] \in Z$  to (S) by the *saddle point residual*

$$\epsilon_{\text{sad}}(z) = \bar{\phi}(z_1) - \underline{\phi}(z_2) = \underbrace{[\bar{\phi}(z_1) - \text{Opt}(P)]}_{\geq 0} + \underbrace{[\text{Opt}(D) - \underline{\phi}(z_2)]}_{\geq 0}. \tag{6}$$

#### 2.1.1 Assumptions and goal

When speaking about a BSP problem (S), our goal is to solve it within a given accuracy  $\epsilon > 0$ , i.e., to find  $z^\epsilon \in Z$  such that  $\epsilon_{\text{sad}}(z^\epsilon) \leq \epsilon$ . Deterministic first order algorithms achieve this goal by working with the values of the associated operator  $F$  at the iterates  $z_t, t = 1, 2, \dots$ , generated by the method. When  $Z$  is simple and the problem is large-scale, computing the values  $F(z_t)$  dominates the computational effort. Our goal in this paper is to replace relatively expensive (in the large scale case)

exact values  $F(z_t)$  with their computationally cheap unbiased random estimates. To this end we assume that

- [P] every point  $z \in Z$  is associated with a probability distribution  $P_z$  such that
  - $P_z$  is supported on  $Z$  and  $\mathbf{E}_{\zeta \sim P_z} \{\zeta\} = z$ ;
  - Given  $z$ , we can sample from the distribution  $P_z$ .

Under these assumptions, and due to the affinity of  $F$ , in order to get an unbiased estimate of  $F(z_t)$ , it suffices to draw a  $\zeta_t \sim P_{z_t}$  and to take  $F(\zeta_t)$  as a desired estimate of  $F(z_t)$ . To make this approach meaningful, the cost of generating  $\zeta_t$  and subsequent computation of  $F(\zeta_t)$  should be significantly less than the cost of a straightforward computation of  $F(z_t)$ . This requirement guided us in the selection of problems to be considered below and in building the s.p. reformulations of these problems.

Note that the deterministic algorithms remain in our scope, since there always is the option to define  $P_z$  as  $\delta_z$  (the unit mass sitting at  $z$ ).

### 2.1.2 Application example: low dimensional approximation

Consider the following problem (related to dimensionality reduction problem in statistics, see, e.g., [4]): given a collection  $V = \{v_1, \dots, v_N\}$  of unit vectors in  $\mathbf{R}^n$ , we want to find a linear subspace  $E \subset \mathbf{R}^n$  of a given dimension  $d < n$  which minimizes the deviation  $\delta(V, E)$  of  $V$  from  $E$ , defined as the worst-case, w.r.t.  $v_i \in V$ , Euclidean distance from  $v_i$  to  $E$ :  $\delta(V, E) = \max_{1 \leq i \leq N} \min_{u \in E} \|v_i - u\|_2$ .

Let  $\Pi^d$  be the family of all orthonormal projectors of  $\mathbf{R}^n$  onto subspaces of dimension  $d$ . Taking into account that  $v_i$  are unit vectors, we have for every  $P \in \Pi^d$ :  $1 - \delta^2(V, \text{Im } P) = \min_i v_i^T P v_i$ , so that

$$\delta_*^2 := \min_E \{\delta^2(V, E) : \dim E = d\} = 1 - \text{Opt}_*, \quad \text{Opt}_* = \max_{P \in \Pi^d} [\min_i v_i^T P v_i].$$

Now, the set  $\Pi^d$  is nonconvex, so that the problem  $\text{Opt}_* = \max_{P \in \Pi^d} [\min_i v_i^T P v_i]$  is seemingly difficult; it, however, admits the tractable relaxation:

$$\text{Opt}_* \leq \text{Opt} := \max_{Q \in \mathcal{P}^d} \min_{1 \leq i \leq N} v_i^T Q v_i, \quad \mathcal{P}^d = \{Q \in \mathbf{S}^n : 0 \leq Q \leq I, \text{Tr}(Q) = d\}. \quad (7)$$

We refer to (7) as the *problem of low dimensional approximation*. We clearly have  $\text{Opt}_* \leq \text{Opt} \leq 1$ , whence  $\delta^2 := 1 - \text{Opt} \leq \delta_*^2$ . Our relaxation admits some quality guarantees. Specifically, let  $Q$  be an optimal solution to (7) and let  $E$  be spanned by the  $d$  leading eigenvectors of  $Q$ . Then

$$\delta(V, E) \leq \delta \sqrt{d+1} \leq \delta_* \sqrt{d+1}. \quad (8)$$

Indeed, let  $e_1, \dots, e_n$  be an orthonormal system of eigenvectors of  $Q$ , and  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  be the corresponding eigenvalues. Note that  $\lambda_k \in [0, 1]$  and  $\lambda_{d+1} \leq \text{Tr}(Q)/(d+1) = d/(d+1)$ . For every  $i$  we have  $\text{Opt} \leq \sum_{k=1}^n \lambda_k (v_i^T e_k)^2 \leq \sum_{k=1}^d (v_i^T e_k)^2 + \lambda_{d+1} \sum_{k=d+1}^n (v_i^T e_k)^2$  and  $\sum_k (v_i^T e_k)^2 = 1$  ( $v_i$  are unit vectors), whence

$$(1 - \lambda_{d+1}) \sum_{k=d+1}^n (v_i^T e_k)^2 \leq 1 - \text{Opt} = \delta^2,$$

that is,  $\sum_{k=d+1}^n (v_i^T e_k)^2 \leq (d + 1)\delta^2$ . Note that the left hand side in this inequality is the squared distance from  $v_i$  to  $E$ , and (8) follows.

Observe that (7) is nothing but the BSP problem:

$$\text{Opt} = \max_{Q \in \mathcal{P}^d} \min_{\lambda \in \Delta_N} \left[ \text{Tr} \left( Q \sum_{i=1}^N \lambda_i v_i v_i^T \right) \right], \quad \Delta_N = \left\{ \lambda \in \mathbf{R}_+^N : \sum_{i=1}^n \lambda_i = 1 \right\}. \tag{9}$$

In terms of  $(S)$ ,  $Z_1 = \Delta_N \subset E_1 = \mathbf{R}^N$ ,  $E_2$  is the space  $\mathbf{S}^n$  of symmetric  $n \times n$  matrices with Frobenius inner product,  $Z_2 = \mathcal{P}^d \subset E_2$ . The associated operator  $F$  is

$$F(z_1, z_2) = F(\lambda, Q) = \left[ \underbrace{[v_1^T Q v_1; \dots; v_N^T Q v_N]}_{F_1(z_2)}; \underbrace{-\sum_{i=1}^N \lambda_i v_i v_i^T}_{F_2(z_1)} \right]. \tag{10}$$

Assuming that the vectors  $v_i$  are dense, the arithmetic cost of computing the value of  $F$  at a given point is  $O(n^2 N)$ . To reduce this cost by randomization, let us specify the distributions  $P_z$  for a given point  $z = (\lambda, Q) \in Z = Z_1 \times Z_2$ . In order to generate  $\zeta \sim P_{(\lambda, Q)}$ , we proceed as follows:

- Given  $\lambda \in \Delta_N$ , we pick  $i \in \{1, \dots, N\}$  at random, with  $\text{Prob}\{i = i\} = \lambda_i$ ,  $1 \leq i \leq N$ , and set  $\zeta_1^i := e_i$ , where  $e_i, i = 1, \dots, N$ , are standard basic orths in  $\mathbf{R}^N$ .
- Given  $Q \in \mathcal{P}^d$ , we build the eigenvalue decomposition  $Q = U \text{Diag}\{q\} U^T$ . Note that  $q \in \Delta_{n,d} := \{q \in \mathbf{R}^n : 0 \leq q_i \leq 1 \forall i, \sum_{i=1}^n q_i = d\}$ . The extreme points of  $\Delta_{n,d}$  are Boolean vectors with exactly  $d$  nonzero entries. There exists a simple algorithm (see Sect. A.1) which, given as input a vector  $q \in \Delta_{n,d}$ , builds in  $O(1) \min\{d, \ln(n)\} n^2$  a.o.  $n$  extreme points  $q^j, 1 \leq j \leq n$ , of  $\Delta_{n,d}$  along with weights  $\mu_j \geq 0, \sum_j \mu_j = 1$ , such that  $q = \sum_j \mu_j q^j$ . We run this algorithm to build  $\{q^j, \mu_j\}_{j=1}^n$ , pick  $J \in \{1, \dots, n\}$  at random, with  $\text{Prob}\{J = j\} = \mu_j, j = 1, \dots, n$ , and set  $\zeta_2^J = U \text{Diag}\{q^J\} U^T$ , which is a projection matrix.
- Finally, we set  $\zeta = [\zeta_1^i; \zeta_2^J] \in \mathcal{P}^d \times \Delta_N$ .

The family of distributions  $P_{(\lambda, Q)}$  clearly satisfies [P]. The ‘‘setup costs’’ for sampling from  $P_{(\lambda, Q)}$  reduce to those of 1) computing the eigenvalue decomposition of  $Q$ , 2) building  $q^1, \dots, q^n, \mu_1, \dots, \mu_n$  (this cost is  $O(n^3 + \min\{d, \ln(n)\} n^2)$  a.o.) and 3) computing the ‘‘cumulative distributions’’  $\{\lambda^i = \sum_{s=1}^i \lambda_s\}_{i=1}^N$  and  $\{\mu^j = \sum_{s=1}^j \mu_s\}_{j=1}^n$  (what amounts to  $O(n + N)$  a.o.). After the setup cost is paid, a sample  $(i, j)$  can be generated at the cost of just  $O(\ln(n + N))$  a.o. Now let us look at the cost of computing  $F(\zeta^{i,j})$  given  $i, j$ . We have

$$F(\zeta^{i,j}) = \left[ \{v_i^T U \text{Diag}\{q^J\} U^T v_i\}_{i=1}^N; -v_i v_i^T \right].$$

Since  $q^j$  has just  $d$  nonzero entries, all equal to 1, let the indices of these entries be  $j_1, \dots, j_d$ , we have  $v_i^T U \text{Diag}\{q^j\} U^T v_i = \sum_{\ell=1}^d (U_{j_\ell}^T v_i)^2$ , where  $U_j$  is  $j$ th column of  $U$ . We see that computing  $F(\zeta^{i,j})$  costs  $O(n^2 + dnN)$  a.o. Thus, the total cost (including that of the setup) of drawing a sample  $\zeta$  from  $P_{(\lambda, Q)}$  and computing  $F(\zeta)$  is

$$O(n^3 + \min\{d, \ln(n)\}n^2 + n^2 + dnN) = O(n^3 + dnN) \text{ a.o.}$$

When  $d \ll n \ll N$ , this cost is much smaller than the cost  $O(n^2N)$  of computing  $F(z)$  at a “general position” point  $z = (\lambda, Q) \in Z$ .

## 2.2 A generalized bilinear saddle point problem

### 2.2.1 The problem

Assume that we are given a *single-parameter family* of bilinear s.p. problems

$$\text{SV}(\rho) = \min_{z_1 \in Z_1} \max_{z_2 \in Z_2} [\phi^\rho(z_1, z_2) := \phi(z_1, z_2) + \rho\psi(z_1, z_2)], \tag{11}$$

where  $\rho \geq 0$  is a parameter and  $\phi(z_1, z_2), \psi(z_1, z_2)$  are bi-affine in  $z_1$  and  $z_2$ . The *generalized bilinear saddle point* (GBSP) problem associated with this family is, by definition, the optimization program

$$\rho_* = \max\{\rho \geq 0 : \text{SV}(\rho) \leq 0\} \tag{12}$$

A highly desirable property of a GBSP problem, relative to our approach, is the convexity of  $\text{SV}(\rho)$  as a function of  $\rho \geq 0$ . To ensure this property, from now on we make the following assumption on the structure of (11):

**[A.1]**  $Z_1 = Z_{11} \times Z_{12}$  is the direct product of two convex compact sets, and the bilinear functions  $\phi(z_1, z_2), \psi(z_1, z_2)$  in (11) are of the form

$$\begin{aligned} \phi(z_1 = [z_{11}; z_{12}], z_2) &= \nu + \langle a_{11}, z_{11} \rangle + \langle b, z_2 \rangle + \langle z_2, Bz_{11} \rangle, \\ \psi(z_1 = [z_{11}; z_{12}], z_2) &= \chi + \langle a_{12}, z_{12} \rangle + \langle c, z_2 \rangle + \langle z_2, Cz_{12} \rangle, \end{aligned} \tag{13}$$

that is,  $\phi(z_1, z_2)$  and  $\psi(z_1, z_2)$  as functions of  $z_1$  depend each on its own “block” of  $z_1$ , and these blocks  $z_{11}$  and  $z_{12}$ , independently of each other, run through the respective convex compact sets  $Z_{11}$  and  $Z_{12}$ .

From now on, we denote by  $F^\rho(z) = \Phi(z) + \rho\Psi(z)$  the affine monotone operator associated with  $\phi^\rho$  according to (4), where  $\Phi(\cdot)$  and  $\Psi(\cdot)$  are the affine monotone operators associated with functions  $\phi(\cdot)$  and  $\psi(\cdot)$ , respectively.

**Lemma 1** *In the case of A.1 the function  $\text{SV}(\rho)$  given by (11) is convex in  $\rho \geq 0$ .*

*Proof* We have

$$\begin{aligned}
 SV(\rho) &= \max_{z_2 \in Z_2} \min_{z_1 \in Z_1} \phi^\rho(z_1, z_2) \\
 &= \max_{z_2 \in Z_2} \min_{z_{11} \in Z_{11}, z_{12} \in Z_{12}} [v + \rho\chi + \langle a_{11}, z_{11} \rangle + \langle b, z_2 \rangle + \langle z_2, Bz_{11} \rangle \\
 &\quad + \rho [\langle a_{12}, z_{12} \rangle + \langle c, z_2 \rangle + \langle z_2, Cz_{12} \rangle]] \\
 &= \max_{z_2 \in Z_2} \left[ v + \rho\chi + \langle b, z_2 \rangle + \rho \langle c, z_2 \rangle \right. \\
 &\quad \left. + \min_{z_{11} \in Z_{11}} \left[ \langle a_{11}, z_{11} \rangle + \langle z_2, Bz_{11} \rangle + \rho \underbrace{\min_{z_{12} \in Z_{12}} [\langle a_{12} + C^* z_2, z_{12} \rangle]}_{g(z_2)} \right] \right] \\
 &= \max_{z_2 \in Z_2} \left[ v + \rho\chi + \langle b, z_2 \rangle + \rho \langle c, z_2 \rangle + \rho g(z_2) \right. \\
 &\quad \left. + \min_{z_{11} \in Z_{11}} \left[ \underbrace{\langle a_{11}, z_{11} \rangle + \langle z_2, Bz_{11} \rangle}_{h(z_2)} \right] \right] \\
 &= \max_{z_2 \in Z_2} \left[ v + \langle b, z_2 \rangle + h(z_2) + \rho [\chi + \langle c, z_2 \rangle + g(z_2)] \right]
 \end{aligned}$$

and thus  $SV(\rho)$  is the supremum of affine functions of  $\rho$ . □

From now on we assume, in addition to **A.1**, that

**[A.2]** *Function  $SV(\rho)$  given by (11) is nonpositive somewhere on  $\mathbf{R}_{++}$  and tends to  $+\infty$  as  $\rho \rightarrow +\infty$ ,*

which implies solvability of (12) and positivity of  $\rho_*$ .

*The goal.* Given a GBSP problem (11)–(12) and a tolerance  $\epsilon > 0$ , our goal will be to find an  $\epsilon$ -solution to the problem, that is, a pair  $\rho_\epsilon, z_1^\epsilon \in Z_1$  such that

$$\rho_\epsilon \geq \rho_* \quad \text{and} \quad \max_{z_2 \in Z_2} \phi^{\rho_\epsilon}(z_1^\epsilon, z_2) \leq \rho_\epsilon \epsilon \tag{14}$$

We are about to point out several important application examples for GBSP problem.

### 2.2.2 Application example: $\ell_1$ minimization with $\ell_p$ fit

Given an  $\ell_p$  norm with  $p \in [1, \infty]$  and a matrix  $A \in \mathbf{R}^{m \times n}$ , the problem of interest is

$$\text{Opt} = \min_x \{ \|x\|_1 : \|Ax - b\|_p \leq \delta \}. \tag{15}$$



Different versions of this problem arise in sparsity-oriented Signal Processing and Compressed Sensing. Setting  $x = \frac{u}{\rho}$ ,  $\|u\|_1 \leq 1$ , we can rewrite the problem in (15) equivalently as

$$1/\text{Opt} = \rho_* = \max \left\{ \rho : \min_{\|u\|_1 \leq 1} \|Au - \rho b\|_p - \rho\delta \leq 0 \right\}. \tag{16}$$

Let  $\varphi(u, \rho) := \|Au - \rho b\|_p - \rho\delta$ . For any given  $u$ , let  $v_u$  be such that  $Au - \rho b = \rho\delta v_u$ . Whenever  $\|v_u\|_p \leq 1$ , we have  $\varphi(u, \rho) = \|Au - \rho b\|_p - \rho\delta = \rho\delta\|v_u\|_p - \rho\delta \leq 0$ . Moreover whenever  $\varphi(u, \rho) \leq 0$ , we see that there exists  $v_u$  satisfying  $\|v_u\|_p \leq 1$  and  $Au - \rho b - \rho\delta v_u = 0$ . Hence we can alternatively write (16) as

$$1/\text{Opt} = \rho_* = \max \left\{ \rho : \Phi(\rho) = \min_{\|u\|_1 \leq 1, \|v\|_p \leq 1} \|Au - \rho b - \rho\delta v\|_\infty \leq 0 \right\}. \tag{17}$$

The advantage of formulation (17) as opposed to (16) (as well as to the original problem given in (15)) lies in the computational complexity of the corresponding first-order oracles. In particular, when computing  $F$  for the BSP's in (17), vector variables participating in nontrivial matrix-vector products vary in unit  $\ell_1$  balls, which, as we shall see, makes an efficient randomization possible. Unfortunately we do not know of extensions of such a randomization for the unit balls of general  $\ell_p$  norms.

Problem given in (17) is nothing but the GBSP problem (11) with  $SV(\rho) = \min_{z_1 \in Z_1} \max_{z_2 \in Z_2} \phi^\rho(z_1, z_2)$ ,

$$\begin{aligned} \phi^\rho(z_1 (= [z_{11}; z_{12}]), z_2) &= z_2^T J_m^T (A J_n z_{11} - \rho[b + \delta z_{12}]), \\ Z_1 &= \underbrace{\Delta_{2n}}_{Z_{11}} \times \underbrace{\{z_{12} \in \mathbf{R}^m : \|z_{12}\|_p \leq 1\}}_{Z_{12}}, \quad Z_2 = \Delta_{2m}, \end{aligned} \tag{18}$$

where we denote  $J_k = [I_k, -I_k]$ ,  $I_k$  being  $k \times k$  identity matrix. This problem satisfies [A.1]; when  $\|b\|_p > \delta$  (otherwise the optimal solution to (15) is  $x = 0$ ), the problem satisfies [A.2] as well. The associated saddle value function is

$$\begin{aligned} SV(\rho) &= \max_{z_2 \in \Delta_{2m}} \min_{z_{11} \in \Delta_{2n}, z_{12} \in Z_{12}} [z_2^T J_m^T (A J_n z_{11} - \rho[b + \delta z_{12}])] \\ &= \max_{w = J_m z_2, z_2 \in \Delta_{2m}} \min_{u = J_n z_1, z_1 \in \Delta_{2n}} \min_{z_{12} \in Z_{12}} [w^T (Au - \rho[b + \delta z_{12}])] \\ &= \max_{\|w\|_1 \leq 1} \min_{\|u\|_1 \leq 1} \min_{\|v\|_p \leq 1} [w^T (Au - \rho[b + \delta v])] = \Phi(\rho). \end{aligned}$$

Suppose that we are given an  $\varepsilon$ -solution  $\rho_\varepsilon$ ,  $z_1^\varepsilon = [z_{11}^\varepsilon; z_{12}^\varepsilon]$  to the problem (14), (18) with  $\varepsilon = \epsilon m^{-\frac{1}{p}}$ . When setting  $x_\varepsilon = \rho_\varepsilon^{-1} J_n z_{11}^\varepsilon$  and  $v_\varepsilon = z_{12}^\varepsilon$  we get an approximate solution to (15) such that

$$\|x_\varepsilon\|_1 \leq \text{Opt} \ \& \ \|Ax_\varepsilon - b\|_p \leq \|\delta v_\varepsilon\|_p + \|Ax_\varepsilon - b - \delta v_\varepsilon\|_p \leq \delta + \epsilon m^{1/p} = \delta + \epsilon.$$

Finally, we associate with  $z = [z_{11}; z_{12}; z_2] \in Z = Z_1 \times Z_2$  a distribution  $P_z$  satisfying condition [P] from Sect. 2.1.1 as follows. Note that for  $z \in Z$ ,  $z_{11}$  and  $z_2$  are vectors from the standard simplices and thus can be considered as probability distributions on the corresponding index sets  $\{1, \dots, 2n\}$ ,  $\{1, \dots, 2m\}$ . To generate

$\zeta = [\zeta_{11}; \zeta_{12}; \zeta_2] \sim P_z$ , we draw at random index  $\iota$  from the distribution  $z_{11}$  and make  $[\zeta_{11}]_{\iota} = 1$  the only nonzero entry in  $\zeta_{11}$ .  $\zeta_2$  is built similarly, with  $z_2$  in the role of  $z_{11}$ , and  $\zeta_{12}$  is nothing but  $z_{12}$ . It is immediately seen that it takes just  $O(m + n)$  a.o. to generate a sample  $\zeta \sim P_z$  and to compute the vector  $F^\rho(\zeta)$ .

It is worth to mention that in the important case  $p = \infty$  the construction of the GBSP which corresponds to (15) can be substantially simplified. Indeed, one can see immediately that for  $p = \infty$  (16) is equivalent to the GBSP problem on the direct product of just two unit  $\ell_1$ -balls (since  $\|Az_1 - b\|_\infty = \max_{\|z_2\|_1 \leq 1} z_2^T (Az_1 - b)$ ). It is more convenient to pass from  $\ell_1$ -balls to the standard simplexes, as it was done in the case of (18). The resulting GBSP problem is given by

$$\begin{aligned} \phi^\rho(z_1, z_2) &= z_2^T J_m^T A J_n z_1 - \rho z_2^T J_m^T b - \rho \delta, \\ Z_1 = Z_{11} &= \Delta_{2n}, \quad Z_{12} = \{0\}, \quad Z_2 = \Delta_{2m}, \end{aligned} \tag{19}$$

and satisfies [A.1] and [A.2] when  $\delta < \|b\|_\infty$ .

### 3 Solving bilinear saddle point problem

We are about to present two randomized first order methods for solving BSPs; they will also be instrumental in solving GBSPs—the *Stochastic Approximation* (SA) and the *Stochastic Mirror Prox* (SMP) algorithms, which are the randomized versions of the methods proposed in [12] and [13], respectively. Both SA and SMP are directly applicable to a BSP problem which we consider in this section; the GBSP case will be considered in Sect. 4.

#### 3.1 The setup

Both SA and SMP algorithms are aimed at solving a BSP problem ( $\mathcal{S}$ ). The setup for these methods is given by

- a norm  $\|\cdot\|$  on the Euclidean space  $E$  where the domain  $Z = Z_1 \times Z_2$  of ( $\mathcal{S}$ ) lives, along with the conjugate norm  $\|\zeta\|_* = \max_{\|z\| \leq 1} \langle \zeta, z \rangle$ ;
- a *distance-generating function* (d.g.f.)  $\omega(z)$  which is convex and continuous on  $Z$ , admits continuous on the set  $Z^\circ = \{z \in Z : \partial\omega(z) \neq \emptyset\}$  selection  $\omega'(z)$  of subgradient (here  $\partial\omega(x)$  is a subdifferential of  $\omega|_Z$  taken at  $z$ ), and is strictly convex with modulus 1 w.r.t.  $\|\cdot\|$ :

$$\forall z', z'' \in Z^\circ : \langle \omega'(z') - \omega'(z''), z' - z'' \rangle \geq \|z' - z''\|^2.$$

We shall refer to the latter property as to *compatibility* of  $\omega(\cdot)$  and  $\|\cdot\|$ .

A d.g.f.  $\omega$  gives rise to several important for us entities:

1. *Bregman distance*  $V_z(u) = \omega(u) - \omega(z) - \langle \omega'(z), u - z \rangle$ , where  $z \in Z^\circ$  and  $u \in Z$ ;
2. *Prox-mapping*  $\text{Prox}_z(\xi) = \text{argmin}_{w \in Z} \{\langle \xi, w \rangle + V_z(w)\} : E \rightarrow Z^\circ$ ; here  $z \in Z^\circ$  is a “prox center;”

3. “ $\omega$ -center”  $z_\omega = \operatorname{argmin}_{z \in Z} \omega(z) \in Z^o$  of  $Z$  and the quantities

$$\Omega = \max_{z \in Z} V_{z_\omega}(z) \leq \max_{z \in Z} \omega(z) - \min_{z \in Z} \omega(z), \quad \Theta = \sqrt{2\Omega}. \tag{20}$$

In the sequel, we set

$$\mathcal{R} := \max_{z \in Z} \|z - z_\omega\| \leq \Theta, \tag{21}$$

where the concluding inequality follows from the fact that for every  $z \in Z$  one has  $\frac{1}{2} \|z - z_\omega\|^2 \leq V_{z_\omega}(z)$  by strong convexity of  $\omega(\cdot)$ .

We also denote by  $\mathcal{L}$  the  $(\|\cdot\|, \|\cdot\|_*)$ -Lipschitz constant of  $F$ :

$$\|F(z) - F(z')\|_* = \|\mathcal{A}(z - z')\|_* \leq \mathcal{L} \|z - z'\|, \quad \forall z, z'; \tag{22}$$

and set

$$M_* = \max_{z, z' \in Z} \|F(z) - F(z')\|_* \leq 2\mathcal{R}\mathcal{L} \leq 2\Theta\mathcal{L}, \tag{23}$$

$$F_* = \max_{z \in Z} \|F(z)\|_* \leq \|F(z_\omega)\|_* + 2\Theta\mathcal{L}. \tag{24}$$

### 3.2 The SA and SMP algorithms

We assume that we have access to an “oracle”  $\mathcal{O}$  which, at  $i$ th call ( $i = 1, 2, \dots$ ), given an input point  $z_i$ , returns a vector  $\xi_i \in E$  such that  $\mathbf{E}_{\xi_i}[\xi_i] = F(z_i)$ . This vector,  $\xi_i$ , can be random with distribution depending on previous calls and, more generally, on the history of our computational process before the call. In fact, in the case when  $\xi_i$  is random, the oracle can be interpreted as providing stochastic subgradient information of the saddle point objective at point  $z_i$ . Whenever this oracle is deterministic, i.e.,  $\xi_i = F(z_i)$ , it is the usual first-order oracle providing subgradient information.

This oracle gives rise to two conceptual algorithms:

$$\begin{aligned} (a) : z_1 = z_\omega; \{z_t, \xi_t\} &\mapsto \{z_{t+1} = \operatorname{Prox}_{z_t}(\gamma_t \xi_t), \xi_{t+1}\}, t = 1, 2, \dots \\ (b) : z_1 = z_\omega; \{z_t, \xi_{2t-1}\} &\mapsto \{w_t = \operatorname{Prox}_{z_t}(\gamma_t \xi_{2t-1}), \xi_{2t}\} \\ &\mapsto \{z_{t+1} = \operatorname{Prox}_{z_t}(\gamma_t \xi_{2t}), \xi_{2t+1}\}, t = 1, 2, \dots \end{aligned} \tag{25}$$

here, in the case of (a),  $z_t$  are the search points, and  $\xi_t$  are the estimates of  $F(z_t)$  as reported by  $\mathcal{O}$ ; in the case of (b),  $z_t, w_t$  are search points, and  $\xi_{2t-1}, \xi_{2t}$  are the estimates of  $F(z_t)$  and  $F(w_t)$ , respectively, as reported by  $\mathcal{O}$ . In both cases,  $\gamma_1, \gamma_2, \dots$  are positive stepsizes defined in a non-anticipative fashion, that is,  $\gamma_t$  depends on oracle’s answers obtained prior to step  $t$  (i.e.,  $\gamma_t$  depends solely on  $\xi_1, \dots, \xi_{t-1}$  in the case of (a), and solely on  $\xi_1, \dots, \xi_{2t-2}$  in the case of (b)). Note that these algorithms (25.a, b) can be perceived as the conceptual versions (with the possibility of working with stochastic oracles) of the Mirror Descent algorithm of [12] and Mirror Prox algorithm of [13], respectively. The main difference of (25.b) from (25.a) is the use of the extra subgradient information. Note that deterministic versions of extra-gradient type algorithms such as Mirror Prox have been shown to be optimal first-order methods (in terms of the number of iterations required for a fixed given accuracy) for structured

non-smooth optimization problems, including s.p. problems over simple domains. For further details on the deterministic versions of these methods, we refer the reader to [12, 13]. Here the main difference of these conceptual algorithms from their deterministic counterparts is as follows: At each iteration, instead of using the exact subgradient information in generation of the next search point, an unbiased random estimate of the current point is built, from which the exact first-order information is gathered and used for computing the next iterate. As opposed to the deterministic versions or the earlier stochastic prototypes of these algorithms, which work with the actual search points, here we average these random estimates to construct the solutions. We refer to (25.a, b) as the stochastic approximation (SA) and stochastic mirror prox (SMP) schemes, respectively. We will consider two implementations of these schemes, the *basic* and the *advanced* ones.

### 3.2.1 Basic implementation

Recall that we have associated with  $(S)$  the affine operator  $F(z) : Z \rightarrow E$  given by (3), and with every point  $z \in Z$ —a probability distribution  $P_z$  supported on  $Z$  satisfying  $\mathbf{E}_{\zeta \sim P_z} \{\zeta\} = z$ . Suppose that

- the stepsizes  $\gamma_t > 0$  are chosen in a non-anticipating fashion such that  $\gamma_1 \geq \gamma_2 \geq \dots$ ;
- in SA:  $\zeta_t$  is drawn at random from the distribution  $P_{z_t}$ , and  $\xi_t = F(\zeta_t)$ ;
- in SMP:  $\xi_{2t-1} = F(\eta_t)$  with  $\eta_t$  drawn at random from the distribution  $P_{z_t}$ , and  $\xi_{2t} = F(\zeta_t)$  with  $\zeta_t$  drawn at random from the distribution  $P_{w_t}$ .

The approximate solution generated by the short-step SA/SMP in course of  $t = 1, 2, \dots$  steps is

$$z^t = t^{-1} \sum_{\tau=1}^t \zeta_\tau. \tag{26}$$

### 3.2.2 Advanced implementation

In Advanced implementation of SA and SMP, same as in the Basic one, the stepsizes  $\gamma_t > 0$  still are chosen in a non-anticipating fashion, but the restriction  $\gamma_1 \geq \gamma_2 \geq \dots$  is now lifted. To explain how the oracle is built, observe that if  $u \in Z$ , then

$$\mathbf{E}_{\zeta \sim P_u} \{\langle F(\zeta), \zeta - u \rangle\} = 0$$

(recall that  $F(z) = a + \mathcal{A}z$  with skew symmetric  $\mathcal{A}$  and that  $\mathbf{E}_{\zeta \sim P_u} \{\zeta\} = u$ ). It follows that given  $u$  and generating one by one independent samples  $\eta^s \sim P_u, s = 1, 2, \dots$ , we will generate with probability 1 a  $\zeta$  such that

$$\langle F(\zeta), \zeta - u \rangle \leq 0. \tag{27}$$

At step  $t$  of SA, in order to define  $\xi_t$ , the oracle draws one by one samples  $\eta^s \sim P_{z_t}, s = 1, 2, \dots$ , until a sample  $\zeta_t := \eta^s$  satisfying (27) with  $u = z_t$  is generated; when it happens, the oracle returns  $\xi_t = F(\zeta_t)$ . At a step  $t$  of SMP, the oracle is invoked twice, first to generate  $\xi_{2t-1} = F(\eta_t)$ , and then to generate  $\xi_{2t} = F(\zeta_t)$ .  $\xi_{2t-1}$  is generated

exactly as in the basic implementation—by drawing a sample  $\eta_t \sim P_{z_t}$  and returning  $\xi_{2t-1} = F(\eta_t)$ . To generate  $\xi_{2t}$ , the oracle draws one by one samples  $\eta^s \sim P_{w_t}$ ,  $s = 1, 2, \dots$ , until a sample  $\zeta_t = \eta^s$  satisfying (27) with  $u = w_t$  is generated; when it happens, the oracle returns  $\xi_{2t} = F(\zeta_t)$ .

Finally, in the advanced implementation we replace the rule (26) for generating approximate solutions with the rule

$$z^t = \frac{1}{\sum_{\tau=1}^t \gamma_\tau} \sum_{\tau=1}^t \gamma_\tau \zeta_\tau. \tag{28}$$

### 3.2.3 Quantifying quality of approximate solutions

Observe that by construction at a step  $\tau$  both  $\zeta_\tau$  and  $F(\zeta_\tau)$  become known. Recalling that  $F$  is affine, it follows that after  $t$  steps we have at our disposal both the approximate solution  $z^t = [z_1^t; z_2^t]$  and the vector  $F(z^t)$ . As a result, with both Basic and Advanced implementations of both SA and SMP, after  $t = 1, 2, \dots$  steps we have at our disposal the quantities

$$\begin{aligned} \bar{\phi}(z_1^t) &= \nu + \langle a_1, z_1^t \rangle + \max_{z_2 \in Z_2} \langle z_2, -F_2(z_1^t) \rangle, \\ \underline{\phi}(z_2^t) &= \nu + \langle a_2, z_2^t \rangle + \min_{z_1 \in Z_1} \langle z_1, F_1(z_2^t) \rangle \end{aligned} \tag{29}$$

[see (3)] and consequently we know the residual  $\epsilon_{\text{sad}}(z^t) = \bar{\phi}(z^t) - \underline{\phi}(z^t)$  of the current approximate solution  $z^t$ . As we shall see in Sect. 4, this feature of our algorithms becomes instrumental when solving GBSP problems.<sup>4</sup> This is in sharp contrast with the prototypes of the SA and the SMP proposed, respectively, in [14, Section 3.3] and [10]. The approximate solutions  $z^t$  of those algorithms were computed according to the formula (28), but with  $z_\tau$  [14] or  $w_\tau$  [10] in the role of  $\zeta_\tau$ . As a result, in the prototype algorithms there is no universal and computationally cheap way to quantify the quality of approximate solutions.

### 3.3 Efficiency estimates for Basic implementation

The accuracy bounds for Basic SA and SMP algorithms are given by the following

**Proposition 2** *Let the BSP problem ( $\mathcal{S}$ ) be solved by the short-step SA or SMP algorithm with positive stepsizes  $\gamma_1 \geq \gamma_2 \geq \dots$  chosen in a non-anticipative fashion. Then*

- (i) *For every  $t \geq 1$ , for both SA and SMP one has*

$$\epsilon_{\text{sad}}(z^t) \leq t^{-1} \left[ \gamma_t^{-1} \Omega + R_t + S_t \right], \quad R_t := \sum_{\tau=1}^t r_\tau, \quad S_t := \sum_{\tau=1}^t s_\tau, \tag{30}$$

<sup>4</sup> Of course, computing the quantities in (29) is not completely costless; note, however, that the cost of one step of the algorithm is dominated by the cost of computing the prox-mapping(s). Thus computing  $\bar{\phi}(\cdot)$  and  $\underline{\phi}(\cdot)$  represents a small fraction of the overall computational effort.

where

$$\begin{aligned}
 r_t &= \begin{cases} \langle F(\zeta_t), \zeta_t - z_t \rangle & \text{in the case of SA,} \\ \langle F(\zeta_t), \zeta_t - w_t \rangle & \text{in the case of SMP,} \end{cases} \\
 s_t &= \begin{cases} \langle F(\zeta_t), z_t - z_{t+1} \rangle - \gamma_t^{-1} V_{z_t}(z_{t+1}), & \text{in the case of SA,} \\ \langle F(\zeta_t), w_t - z_{t+1} \rangle - \gamma_t^{-1} V_{z_t}(z_{t+1}), & \text{in the case of SMP.} \end{cases}
 \end{aligned}$$

We have

$$s_t \leq \begin{cases} \frac{\gamma_t}{2} \|F(\zeta_t)\|_*^2, & \text{in the case of SA,} \\ \frac{\gamma_t}{2} \|F(\zeta_t) - F(\eta_t)\|_*^2 - \frac{1}{2\gamma_t} \|w_t - z_t\|^2, & \text{in the case of SMP,} \end{cases} \tag{31}$$

implying

$$s_t \leq \begin{cases} \frac{\gamma_t}{2} F_*^2, & \text{in the case of SA,} \\ \frac{\gamma_t}{2} M_*^2, & \text{in the case of SMP.} \end{cases} \tag{32}$$

In particular, if the stepsizes  $\gamma_t > 0$  satisfy  $S_t \leq \Omega/\gamma_t$ ,  $t = 1, 2, \dots$ , then

$$\epsilon_{\text{sad}}(z^t) \leq \frac{2\Omega}{t\gamma_t} + \frac{R_t}{t}. \tag{33}$$

(ii) Further,  $\mathbf{E}\{R_t\} = 0$ , and in the case of SMP, under additional assumption that

$$\gamma_t \leq (\sqrt{3}\mathcal{L})^{-1}, \tag{34}$$

we have

$$s_t \leq \frac{3\gamma_t}{2} \left[ \|\mathcal{A}(\zeta_t - w_t)\|_*^2 + \|\mathcal{A}(\eta_t - z_t)\|_*^2 \right], \tag{35}$$

so that  $\mathbf{E}\{s_t\} \leq 3\gamma_t\sigma^2$ , where

$$\sigma^2 = \sup_{z \in Z} \mathbf{E}_{\zeta \sim P_z} \left\{ \|\mathcal{A}(\zeta - z)\|_*^2 \right\} \leq M_*^2. \tag{36}$$

In particular, if the stepsizes  $\gamma_t > 0$  satisfy  $\mathbf{E}\{S_t\} \leq \Omega/\gamma_t$  for  $t = 1, 2, \dots$ , then

$$\mathbf{E}\{\epsilon_{\text{sad}}(z^t)\} \leq \frac{2\Omega}{t\gamma_t}.$$

*Proof 1<sup>0</sup>.* We need the following

**Lemma 2** [cf. [13], Lemma 3.1.(b)] *Given  $z \in Z^o$ ,  $\gamma > 0$  and  $\xi, \eta \in E$ , let us set*

$$w = \text{Prox}_z(\gamma\xi) = \underset{v \in Z}{\operatorname{argmin}} \{ \langle \gamma\xi - \omega'(z), v \rangle + \omega(v) \},$$

$$z_+ = \text{Prox}_z(\gamma\eta) = \underset{v \in Z}{\operatorname{argmin}} \{ \langle \gamma\eta - \omega'(z), v \rangle + \omega(v) \}.$$

Then  $w, z_+ \in Z^o$ , and for every  $u \in Z$  one has

$$\begin{aligned} (a) \quad & \gamma \langle \eta, w - u \rangle \leq V_z(u) - V_{z_+}(u) + \gamma \langle \eta, w - z_+ \rangle - V_z(z_+) \\ (b) \quad & \leq V_z(u) - V_{z_+}(u) + \gamma \langle \eta - \xi, w - z_+ \rangle - V_z(w) - V_w(z_+) \\ (c) \quad & \leq V_z(u) - V_{z_+}(u) + \gamma \|\eta - \xi\|_* \|w - z_+\| - \frac{1}{2} \left[ \|w - z\|^2 + \|w - z_+\|^2 \right] \\ (d) \quad & \leq V_z(u) - V_{z_+}(u) + \frac{1}{2} \left[ \gamma^2 \|\eta - \xi\|_*^2 - \|w - z\|^2 \right]. \end{aligned} \tag{37}$$

*Proof of Lemma 2* The inclusions  $w, z_+ \in Z^o$  are evident (a subgradient of  $\omega(\cdot)$  at  $w$ , taken w.r.t.  $Z$ , is, e.g.,  $\omega'(z) - \gamma\xi$ , and similarly for  $z_+$ ). Now let  $u \in Z$ .  $z_+$  is an optimal solution of certain explicit convex optimization problem; taking into account that  $\omega'(\cdot)$  is continuous on  $Z^o$ , it is easily seen that the necessary optimality condition in this problem reads  $\langle \gamma\eta + \omega'(z_+) - \omega'(z), u - z_+ \rangle \geq 0$ , whence  $\gamma \langle \eta, w - u \rangle \leq \gamma \langle \eta, w - z_+ \rangle + \langle \omega'(z_+) - \omega'(z), u - z_+ \rangle$ , and the latter inequality, after rearranging terms in the right hand side, becomes (a). By similar reasons,  $0 \leq \langle \gamma\xi + \omega'(w) - \omega'(z), v - w \rangle$  for all  $v \in Z$ ; setting  $v = z_+$ , summing up the resulting inequality with (a) and rearranging terms in the right hand side of what we get, we arrive at (b). (c) follows from (b) due to  $V_a(b) \geq \frac{1}{2} \|a - b\|^2$  (recall that  $\omega$  is strongly convex, modulus 1 w.r.t.  $\|\cdot\|$ , on  $Z$ ). Finally, (d) follows from (c) due to  $\mu v - \frac{1}{2} \mu^2 \leq \frac{1}{2} v^2$ .  $\square$

**2<sup>0</sup>.** Let us prove the bound (30). Consider first the case of SMP. Applying Lemma 2 to  $z = z_\tau$ ,  $\gamma = \gamma_\tau$ ,  $\xi = F(\eta_\tau)$ ,  $\eta = F(\zeta_\tau)$ , which results in  $w = w_\tau$  and  $z_+ = z_{\tau+1}$ , we get for all  $u \in Z$ :

$$\gamma_\tau \langle F(\zeta_\tau), w_\tau - u \rangle \leq V_{z_\tau}(u) - V_{z_{\tau+1}}(u) + [\gamma_\tau \langle F(\zeta_\tau), w_\tau - z_{\tau+1} \rangle - V_{z_\tau}(z_{\tau+1})]$$

whence for all  $u \in Z$

$$\begin{aligned} \langle F(\zeta_\tau), \zeta_\tau - u \rangle & \leq \gamma_\tau^{-1} (V_{z_\tau}(u) - V_{z_{\tau+1}}(u)) + \overbrace{\langle F(\zeta_\tau), \zeta_\tau - w_\tau \rangle}^{r_\tau} + s_\tau, \\ s_\tau & = \langle F(\zeta_\tau), w_\tau - z_{\tau+1} \rangle - \gamma_\tau^{-1} V_{z_\tau}(z_{\tau+1}) \\ & \leq \frac{1}{2} \left[ \gamma_\tau \|F(\zeta_\tau) - F(\eta_\tau)\|_*^2 - \gamma_\tau^{-1} \|w_\tau - z_\tau\|^2 \right] \end{aligned} \tag{38}$$

with (\*) given by (37). When summing up inequalities (38) over  $\tau$  and taking into account that  $\gamma_1 \geq \gamma_2 \geq \dots$ ,  $V_z(u) \geq 0$  and  $V_{z_1}(u) = V_{z_\omega}(u) \leq \Omega$  by definition of  $\Omega$ ,

we get

$$\sum_{\tau=1}^t \langle F(\zeta_\tau), \zeta_\tau - u \rangle \leq \gamma_t^{-1} \Omega + \sum_{\tau=1}^t [s_\tau + r_\tau]. \tag{39}$$

On the other hand, taking into account that  $\mathcal{A}$  is skew symmetric,

$$\begin{aligned} \sum_{\tau=1}^t \langle F(\zeta_\tau), \zeta_\tau - u \rangle &= \sum_{\tau=1}^t \langle a + \mathcal{A}\zeta_\tau, \zeta_\tau - u \rangle \\ &= t \langle a, z^t - u \rangle - \sum_{\tau=1}^t \langle \mathcal{A}\zeta_\tau, u \rangle \\ &= t [\langle a, z^t - u \rangle - \langle \mathcal{A}z^t, u \rangle] \\ &= t [\langle a, z^t - u \rangle + \langle \mathcal{A}z^t, z^t - u \rangle] \\ &= t \langle F(z^t), z^t - u \rangle. \end{aligned}$$

Thus, for all  $u \in Z$  it holds

$$t \langle F(z^t), z^t - u \rangle \leq \Omega \gamma_t^{-1} + \sum_{\tau=1}^t [s_\tau + r_\tau] = \gamma_t^{-1} \Omega + S_t + R_t. \tag{40}$$

Setting  $z^t = [z_1^t; z_2^t]$  and  $u = [u_1; u_2]$ , we get from the definition of  $F(\cdot)$  and the bilinearity of the inner product  $\langle F(z^t), z^t - u \rangle = \phi(z_1^t, u_2) - \phi(u_1, z_2^t)$ ; the supremum of the latter quantity over  $u \in Z$  is the s.p. residual  $\epsilon_{\text{sad}}(z^t)$ . Since the right hand side in (40) is independent of  $u$ , we arrive at the SMP-version of (30).

**3<sup>0</sup>.** Now consider the case of SA. Applying Lemma 2 to  $\gamma = \gamma_\tau, z = z_\tau, \xi = 0, \eta = F(\zeta_\tau)$ , which results in  $w = z_\tau$  and  $z_+ = z_{\tau+1}$ , and acting exactly as in the case of SMP, we arrive at the SA-version of (30).

**4<sup>0</sup>.** Let us prove (ii). The conditional to the ‘‘past’’ (the answers of the oracle prior to the call for  $\xi_{2\tau}$ ) distribution of  $\zeta_\tau$  is  $P_{w_\tau}$ , which combines with the affinity of  $F$  and the facts that the linear part of  $F$  is skew symmetric and the expectation of  $P_z$  is  $z$ , to imply that

$$\begin{aligned} \mathbf{E}\{\langle F(\zeta_\tau), \zeta_\tau - w_\tau \rangle\} &= \langle a, \mathbf{E}\{\zeta_\tau\} - w_\tau \rangle + \mathbf{E}\{\langle \mathcal{A}\zeta_\tau, \zeta_\tau - w_\tau \rangle\} = -\mathbf{E}\{\langle \mathcal{A}\zeta_\tau, w_\tau \rangle\} \\ &= \mathbf{E}\{\langle \mathcal{A}(w_\tau - \zeta_\tau), w_\tau \rangle\} = 0, \end{aligned}$$

whence  $\mathbf{E}\{R_t\} = 0$  for all  $t$ . By completely similar reasoning,  $\mathbf{E}\{R_t\} = 0$  in the case of SA. To complete the proof (ii), we need to prove (35). We have

$$\begin{aligned} s_t &\leq \frac{\gamma_t}{2} \|F(\zeta_t) - F(\eta_t)\|_*^2 - \frac{1}{2\gamma_t} \|w_t - z_t\|^2 \quad [\text{see (31)}] \\ &\leq \frac{\gamma_t}{2} [\|F(w_t) - F(z_t)\|_* + \|F(\zeta_t) - F(w_t)\|_* + \|F(\eta_t) - F(z_t)\|_*]^2 \end{aligned}$$



$$\begin{aligned}
 & -\frac{1}{2\gamma_t} \|w_t - z_t\|^2 \\
 \leq & \underbrace{\left[ \frac{3\gamma_t}{2} \mathcal{L}^2 - \frac{1}{2\gamma_t} \right]}_{\leq 0 \text{ by (34)}} \|w_t - z_t\|^2 + \frac{3\gamma_t}{2} \left[ \|F(\zeta_t) - F(w_t)\|_*^2 + \|F(\eta_t) - F(z_t)\|_*^2 \right].
 \end{aligned}$$

It remains to note that  $\|F(\zeta_t) - F(w_t)\|_*^2 + \|F(\eta_t) - F(z_t)\|_*^2 \leq 2M_*^2$  since  $\zeta_t, w_t, \eta_t, z_t \in Z$  and that the expectations of  $\|F(\zeta_t) - F(w_t)\|_*^2$  and  $\|F(\eta_t) - F(z_t)\|_*^2$ , conditional over the respective pasts, do not exceed  $\sigma^2$ .  $\square$

The bound of Proposition 2 allows to easily conceive stepsize policies. Let us start with *offline* policies, where  $\gamma_t$  are chosen in advance deterministic reals. If the number of steps  $N$  is fixed in advance, one can use constant stepsizes  $\gamma_1 = \dots = \gamma_N = \gamma$ . In particular, when choosing

$$\gamma = \begin{cases} \frac{1}{F_*} \sqrt{\frac{2\Omega}{N}}, & \text{in the case of SA (a)} \\ \min \left\{ \frac{1}{\sigma} \sqrt{\frac{\Omega}{3N}}, \frac{1}{\sqrt{3\mathcal{L}}} \right\}, & \text{in the case of SMP (b)} \end{cases} \tag{41}$$

[by (32), (24) this choice implies that  $\mathbf{E}\{S_t\} \leq \Omega/\gamma_t, 1 \leq t \leq N$ ], Proposition 2 implies the efficiency bound

$$\mathbf{E}\{\epsilon_{\text{sad}}(z^N)\} \leq \begin{cases} F_* \sqrt{\frac{2\Omega}{N}}, & \text{in the case of SA (a)} \\ \max \left\{ 2\sigma \sqrt{\frac{3\Omega}{N}}, \frac{2\sqrt{3}\Omega\mathcal{L}}{N} \right\}, & \text{in the case of SMP (b)} \end{cases} \tag{42}$$

When the number of steps is not fixed in advance, one can use the decreasing stepsizes

$$\forall t \geq 1, \gamma_t = \begin{cases} \frac{1}{F_*} \sqrt{\frac{\Omega}{t}}, & \text{in the case of SA,} \\ \min \left\{ \frac{1}{\sigma} \sqrt{\frac{\Omega}{6t}}, \frac{1}{\sqrt{3\mathcal{L}}} \right\}, & \text{in the case of SMP,} \end{cases} \tag{43}$$

which result in the accuracy bound

$$\forall t \geq 1, \mathbf{E}\{\epsilon_{\text{sad}}(z^t)\} \leq \begin{cases} 2F_* \sqrt{\frac{\Omega}{t}}, & \text{in the case of SA (a)} \\ \max \left\{ 2\sigma \sqrt{\frac{6\Omega}{t}}, \frac{2\sqrt{3}\Omega\mathcal{L}}{t} \right\}, & \text{in the case of SMP (b)} \end{cases} \tag{44}$$

completely similar to (42).

### 3.3.1 Online stepsize policies

From theoretical viewpoint, the main advantage of the offline stepsize policies (41) and (43) is that in the framework of our approach they result in the best possible (and

in fact—the best known under circumstances) efficiency estimates (42), (44). While they may appear attractive also from the practical viewpoint because of their apparent simplicity, their use may present several disadvantages: the quantity  $\sigma$  involved in the stepsize computation may not be available at hand and should be evaluated. Besides this, these policies are offline and worst-case oriented; we would prefer more flexible online adjustable stepsizes.

A natural way to adjust the stepsizes online would be to choose at each step  $t \geq 1$  the largest  $\gamma_t \leq \gamma_{t-1}$  ensuring the balance  $\Omega/\gamma_t \geq S_t$ , and thus the bound (33). This idea cannot be implemented “as is”, since the stepsize policy should be non-anticipative, while  $s_t$  is not yet available when  $\gamma_t$  is computed. This difficulty can be easily circumvented by using instead of  $s_t$  its a priori upper bound, which is either  $\frac{\gamma_t}{2} F_*$  for the SA algorithm or  $\frac{\gamma_t}{2} M_*^2$  for the SMP, see (31). Specifically, consider the online policy of choosing  $\gamma_t, t \geq 1$  as follows:

$$\Omega\gamma_t^{-2} = \begin{cases} 2\sum_{\tau=1}^{t-1}\gamma_\tau^{-1}[s_\tau]_+ + F_*^2 & \text{in the case of SA,} \\ 2\sum_{\tau=1}^{t-1}\gamma_\tau^{-1}[s_\tau]_+ + 8\Omega\mathcal{L}^2 & \text{in the case of SMP,} \end{cases} \tag{45}$$

where we set  $\sum_{\tau=1}^0\gamma_\tau^{-1}[s_\tau]_+ = 0$ . With this policy, one clearly has  $\gamma_1 \geq \gamma_2 \geq \dots$ .

**Proposition 3** *Let positive stepsizes  $\gamma_t, t = 1, 2, \dots$  of the Basic SA/SMP implementation be chosen according to (45). Then the approximate solution  $z^t$  satisfies*

$$\epsilon_{\text{sad}}(z^t) \leq \frac{(1 + \sqrt{2})\Omega}{t\gamma_t} + \frac{R_t}{t}. \tag{46}$$

As a consequence, we have

$$\epsilon_{\text{sad}}(z^t) \leq \begin{cases} \left[ \frac{(1+\sqrt{2})\sqrt{\Omega}}{t} \left( F_*^2 + \sum_{\tau=1}^{t-1} \|F(\zeta_\tau)\|_*^2 \right)^{1/2} + \frac{R_t}{t} \right], & \text{in the case of SA (a)} \\ \left[ \frac{(1+\sqrt{2})\sqrt{\Omega}}{t} \left( 8\Omega\mathcal{L}^2 + \sum_{\tau=1}^{t-1} \varsigma_\tau \right)^{1/2} + \frac{R_t}{t} \right] \\ \leq \frac{7\Omega\mathcal{L}}{t} + \frac{R_t}{t} + \frac{(1+\sqrt{2})\sqrt{\Omega}}{t} \sqrt{\sum_{\tau=1}^{t-1} \varsigma_\tau}, & \text{in the case of SMP (b)} \end{cases} \tag{47}$$

where

$$\varsigma_t = 3 \left[ \|F(\zeta_t) - F(w_t)\|_*^2 + \|F(\eta_t) - F(z_t)\|_*^2 \right]. \tag{48}$$

Recalling that  $\mathbf{E}\{R_t\} = 0$  and  $\mathbf{E}\{\varsigma_t\} \leq 6\sigma^2$  [see (36)], we arrive at

**Corollary 1** *Under the premise of Proposition 3, for the SMP algorithm one has*

$$\mathbf{E}\{\epsilon_{\text{sad}}(z^t)\} \leq \frac{7\Omega\mathcal{L}}{t} + \frac{6\sqrt{\Omega}\sigma}{\sqrt{t}}. \tag{49}$$

*Proof of Proposition 3* Observe that (45) implies that  $\gamma_1 \geq \gamma_2 \geq \dots$ .

**1<sup>0</sup>.** Let us verify first that with the choice (45) of  $\gamma_\tau$ ,  $\tau = 1, 2, \dots$  we have for all  $t = 1, 2, \dots$ ,

$$\sqrt{2}\Omega\gamma_t^{-1} \geq S_t. \tag{50}$$

Indeed, for  $t = 2, 3, \dots$  we have (with  $2S_0 = F_*^2$  in the case of SA and  $2S_0 = 8\Omega\mathcal{L}^2, \geq M_*^2$  by (23), in the case of SMP)

$$\frac{\gamma_{t-1}^2}{\gamma_t^2} = \frac{\sum_{\tau=1}^{t-1} 2[s_\tau]_+/\gamma_\tau + 2S_0}{\sum_{\tau=1}^{t-2} 2[s_\tau]_+/\gamma_\tau + 2S_0} \leq 1 + \frac{2[s_{t-1}]_+/\gamma_{t-1}}{2S_0} \leq 2 \tag{51}$$

[recall that  $2s_t/\gamma_t \leq 2S_0$  by (32)]. On the other hand

$$\gamma_t^{-2} - \gamma_{t-1}^{-2} = \frac{2[s_{t-1}]_+}{\Omega\gamma_{t-1}}, \tag{52}$$

and

$$\begin{aligned} \gamma_t^{-1} - \gamma_{t-1}^{-1} &\geq \frac{\gamma_t}{2}(\gamma_t^{-2} - \gamma_{t-1}^{-2}) = \frac{\gamma_t[s_{t-1}]_+}{\gamma_{t-1}\Omega} \geq \frac{[s_{t-1}]_+}{\sqrt{2}\Omega} \Rightarrow \\ \sqrt{2}\Omega[\gamma_t^{-1} - \gamma_{t-1}^{-1}] &\geq [s_{t-1}]_+ \end{aligned} \tag{53}$$

where the second inequality follows from  $\gamma_{t-1} \leq \sqrt{2}\gamma_t$  as implied by (51). By summing up the resulting inequalities in (53), we get

$$\sqrt{2}\Omega\gamma_t^{-1} \geq \sum_{\tau=1}^{t-1} s_\tau + \sqrt{2}\Omega\gamma_1^{-1}. \tag{54}$$

In the case of SMP, we have  $\gamma_1 = (2\sqrt{2}\mathcal{L})^{-1}$ , whence  $\sqrt{2}\Omega\gamma_1^{-1} = 4\Omega\mathcal{L} \geq \gamma_1 M_*^2 \geq \gamma_t M_*^2$  [see (23)], whence  $\sqrt{2}\Omega\gamma_1^{-1} \geq s_t$  in view of (31), and (54) implies (50). In the case of SA, we have  $\gamma_1 = \sqrt{\Omega}/F_*$ , whence  $\sqrt{2}\Omega\gamma_1^{-1} = \sqrt{2}\sqrt{\Omega}F_* \geq \gamma_1 F_*^2 \geq \gamma_t F_*^2$ , whence  $\sqrt{2}\Omega\gamma_1^{-1} \geq s_t$  by (31), and (50) again is given by (54).

**2<sup>0</sup>.** Invoking (30), (50) implies (46). Now, by (31) in the case of SA we have  $2[s_\tau]_+/\gamma_\tau \leq \|F(\zeta_\tau) - F(\eta_\tau)\|_*^2$ . In the case of SMP we have

$$\begin{aligned} \frac{2[s_\tau]_+}{\gamma_\tau} &\leq \|F(\zeta_\tau) - F(\eta_\tau)\|_*^2 - \gamma_\tau^{-2}\|w_\tau - z_\tau\|^2 \quad [\text{see (31)}] \\ &\leq [\|F(\zeta_\tau) - F(w_\tau)\|_* + \|F(w_\tau) - F(z_\tau)\|_* + \|F(z_\tau) - F(\eta_\tau)\|_*]^2 - \gamma_\tau^{-2}\|w_\tau - z_\tau\|^2 \\ &\leq 3[\|F(\zeta_\tau) - F(w_\tau)\|_*^2 + \|F(z_\tau) - F(\eta_\tau)\|_*^2] + [3\|F(w_\tau) - F(z_\tau)\|_*^2 - \gamma_\tau^{-2}\|w_\tau - z_\tau\|^2] \\ &\leq 3[\underbrace{\|F(\zeta_\tau) - F(w_\tau)\|_*^2 + \|F(z_\tau) - F(\eta_\tau)\|_*^2}_{:=\zeta_\tau}] \quad [\text{by (22) due to } \gamma_1^{-1} = 2\sqrt{2}\mathcal{L}] \end{aligned}$$

Invoking (45), we get

$$\gamma_t^{-1} \leq \Omega^{-1/2} \cdot \begin{cases} \left( F_*^2 + \sum_{\tau=1}^{t-1} \|F(\zeta_\tau)\|_*^2 \right)^{1/2}, & \text{in the case of SA} \\ \left( 8\Omega\mathcal{L}^2 + \sum_{\tau=1}^{t-1} \varsigma_\tau \right)^{1/2}, & \text{in the case of SMP} \end{cases} \tag{55}$$

which combines with (46) to imply (47). □

Note that the bounds (47.a) and (49) within an absolute constant factor coincide with the respective bounds in (44), that is, our online stepsizes policy (which, in contrast to (43), does not require knowledge of  $\sigma$ ) is not worse than the “theoretically optimal” stepsize policies underlying (44).

### 3.3.2 Discussion

By definitions of  $F_*$  and  $\sigma$  we have  $\sigma \leq 2F_*$  [see (24), (36)]. As a result, the SA efficiency estimate (44.a) for large  $t$  (and even for all  $t$ , provided that  $F_*$  is of order of  $\sqrt{\Theta\mathcal{L}}$ ) can be better than the SMP efficiency estimate (44.b) by *at most* an absolute constant factor, and becomes much worse than the SMP estimate when  $t$  is large and  $\sigma \ll \sqrt{\Omega\mathcal{L}}$ . Besides this, when the noise level  $\sigma$  of the oracle is small enough (specifically,  $\sigma^2 = O\left(\frac{\Omega\mathcal{L}^2}{N}\right)$ ), the efficiency estimate of SMP satisfies  $\mathbf{E}\{\epsilon_{\text{sad}}(z^t)\} \leq O(1)\frac{\Omega\mathcal{L}}{N}$ , which, modulo expectation of the residual instead of the residual itself, coincides with the best known so far efficiency estimate of the deterministic first order algorithms for solving BSP problems. In addition to this, we do have a possibility to make  $\sigma$  small. The trivial way to do so is to use  $P_z = \delta_z$ , which results in  $\sigma = 0$  and makes SMP a version of the deterministic mirror prox algorithm (DMP) proposed in [13]. Another, more attractive, option to control  $\sigma$  is as follows. Given the family of distributions  $P_z$  supported on  $Z$  and such that  $\mathbf{E}_{\zeta \sim P_z}\{\zeta\} = z$ , and a positive integer  $k$ , we can convert  $P_z$  into the family of distributions  $P_z^{(k)}$ , also supported on  $Z$  and satisfying  $\mathbf{E}_{\zeta \sim P_z}\{\zeta\} = z$ , as follows. In order to generate a random vector  $\zeta \sim P_z^{(k)}$  and to compute  $F(\zeta)$ , we draw a  $k$ -element sample  $\zeta^1, \dots, \zeta^k$  from the distribution  $P_z$ , compute  $F(\zeta^1), \dots, F(\zeta^k)$  and then set  $\zeta = \frac{1}{k}\sum_{i=1}^k \zeta^i$ , so that  $F(\zeta) = \frac{1}{k}\sum_{i=1}^k F(\zeta^i)$ . If, as in the examples of Sect. 2, drawing  $\zeta^i \sim P_z$  and computing  $F(\zeta^i)$  is much cheaper than computing  $F(z)$ , the outlined procedure with a “reasonably large” value of  $k$  is still significantly cheaper than the direct computation of  $F(z)$ . At the same time, for “good enough” norms  $\|\cdot\|_*$ , passing from  $P_z$  to  $P_z^{(k)}$  can significantly reduce the noise level  $\sigma$ . Specifically, given a norm  $\|\cdot\|_*$  on a finite-dimensional Euclidean space  $E$ , one can associate with it its *regularity parameter*  $\varkappa \geq 1$  (see Sect. A.2) to ensure the following: whenever  $k > 0$  is an integer and  $\xi^1, \dots, \xi^k$  are independent vectors from  $E$  with  $\mathbf{E}\{\xi^i\} = 0$  and  $\mathbf{E}\{\|\xi^i\|_*^2\} \leq \alpha_i^2$  and  $\alpha = \max_i \alpha_i$ , then for  $\xi = \frac{1}{k}\sum_{i=1}^k \xi^i$  the following holds

$$\mathbf{E}\{\|\xi\|_*^2\} \leq \min\left[\frac{1}{k}, \frac{\varkappa}{k^2}\right] \sum_{i=1}^k \alpha_i^2 \leq \min\left[1, \frac{\varkappa}{k}\right] \alpha^2.$$

Suppose now that when running SMP we sample  $\zeta_t, \eta_t$  from the distributions  $P_z^{(k)}$  for some  $k > 0$ . It follows that if  $\|\cdot\|_*$  is  $\kappa$ -regular with certain  $\kappa$ , then, passing from  $P_z$  to  $P_z^+ = P_z^{(k)}$ , we can reduce the “original” value of  $\sigma$  to the value  $\sigma^+ = \min[1, \sqrt{\frac{\kappa}{k}}]\sigma$ . We shall see in a while that in the applications we have mentioned so far,  $\kappa$  is “small”—at most logarithmic in  $\dim Z$ . The bottom line is that there is a tradeoff between the computational cost of a call to a stochastic oracle and the noise level  $\sigma$ . Consequently, in the case of SMP, it is possible to tradeoff the computational effort per iteration and the iteration count to obtain an approximate solution of the desired expected quality, and we can use this tradeoff in order to save on the overall amount of computations. This option (which is the major advantage of SMP as compared to SA) is especially attractive when among the two components of our computational effort per iteration—one related to computing  $\eta_t, \zeta_t, F(\eta_t) F(\zeta_t)$ , and the other aimed at computing the prox mappings—the second component is essentially more significant than the first one. In such a situation, we basically can only gain by passing from  $P_z$  to  $P_z^{(k)}$  with  $k$  chosen to balance the outlined two components of the computational effort.

### 3.3.3 Large deviations

In the above efficiency estimates, say, in (49), we upper-bounded the *expected* inaccuracy of approximate solutions  $z^t$ . In fact, one can get exponential upper bounds on probabilities of large deviations for the inaccuracy of the approximate solution. Though we do not need such bounds to access the inaccuracy of solutions, they are still useful to provide *theoretical* guarantees for the complexity of our algorithms (cf. Theorem 1 in the next section).

For the sake of definiteness, when presenting large deviation results, we restrict ourselves to the SMP algorithm and the stepsize strategy (45). We can easily bound from above the probability of  $\epsilon_{\text{sad}}(z^t)$  to be larger than the bound (49) on its expectation using the Markov inequality. Moreover, let us fix the number  $t$  of iterations, run the algorithm  $m$  times and select the best, in terms of  $\epsilon_{\text{sad}}(\cdot)$ , of the resulting approximate solutions. The probability that for this solution  $\epsilon_{\text{sad}}(\cdot)$  is worse than, say, twice the right hand side of (49) is at most  $2^{-m}$  and thus can be made negligibly small for quite moderate values of  $m$ .

We also have the following bound on the deviations of the algorithm without restarts:

**Proposition 4** *Assume we are solving problem (S) by Basic implementation of SMP where  $\zeta_t, \eta_t$  are sampled from the distributions  $P_z^{(k)}$ ,  $k \geq 1$  being a parameter of the construction. Assume also that the norm  $\|\cdot\|_*$  is  $\kappa$ -regular, and the online stepsize policy (45) is used. Then there are absolute constants  $K_0, K_1$  such that the approximate solution  $z^t$  satisfies for all  $t \geq 1$  and  $\lambda, \Lambda \geq 0$*

$$\text{Prob} \left\{ \epsilon_{\text{sad}}(z^t) \geq K_0 \left[ \frac{\Theta^2 \mathcal{L}}{t} + \frac{\kappa_*(k, \Lambda) \Theta^2 \mathcal{L}}{\sqrt{kt}} + \Theta_{F_*} \sqrt{\frac{\lambda}{kt}} \right] \right\} \leq e^{-\Lambda t} + e^{-\lambda}, \quad (56)$$

where  $\kappa_*(k, \Lambda) = \sqrt{\min[k, (\kappa + \Lambda)]}$ . In particular, one has for all  $\varepsilon > 0$ :

$$\begin{aligned} \text{Prob}\{\epsilon_{\text{sad}}(z^N) \geq \varepsilon\} &\leq e^{-\Lambda N} + e^{-\lambda} \text{ for } N \geq N_\varepsilon, \text{ where} \\ N_\varepsilon &= K_1 \text{Ceil} \left( \max \left[ \Theta^2 \mathcal{L} \varepsilon^{-1}, \frac{\kappa_*^2(k, \Lambda) \Theta^4 \mathcal{L}^2}{k \varepsilon^2}, \frac{F_*^2 \Theta^2 \lambda}{k \varepsilon^2} \right] \right). \end{aligned} \tag{57}$$

For proof, see Sect. A.3.

### 3.4 Efficiency estimates for advanced implementation

The efficiency of Advanced implementations of SA and SMP stem from the following result (we use the notation from Sect. 3.1):

**Proposition 5** *Let the BSP problem (S) be solved by the advanced-step SA or SMP algorithms. Then for every  $t \geq 1$ , for both SA and SMP one has*

$$\epsilon_{\text{sad}}(z^t) \leq \Gamma_t^{-1} [\Omega + R_t + S_t] = \Gamma_t^{-1} \left[ \Omega + \sum_{\tau=1}^t r_\tau + \sum_{\tau=1}^t s_\tau \right], \tag{58}$$

where

$$\begin{aligned} \Gamma_t &= \sum_{\tau=1}^t \gamma_\tau, \\ r_t &= \begin{cases} \gamma_t \langle F(\zeta_t), \zeta_t - z_t \rangle & \text{in the case of SA} \\ \gamma_t \langle F(\zeta_t), \zeta_t - w_t \rangle & \text{in the case of SMP} \end{cases} \\ s_t &= \begin{cases} [\gamma_t \langle F(\zeta_t), z_t - z_{t+1} \rangle - V_{z_t}(z_{t+1})], & \text{in the case of SA} \\ [\gamma_t \langle F(\zeta_t), w_t - z_{t+1} \rangle - V_{z_t}(z_{t+1})], & \text{in the case of SMP} \end{cases} \end{aligned}$$

with  $r_t \leq 0$  and

$$s_t \leq \begin{cases} \frac{\gamma_t^2}{2} \|F(\zeta_t)\|_*^2 \leq \frac{\gamma_t^2}{2} F_*^2, & \text{in the case of SA} \\ \frac{\gamma_t^2}{2} \|F(\zeta_t) - F(\eta_t)\|_*^2 - \frac{1}{2} \|w_t - z_t\|^2 \leq \frac{\gamma_t^2}{2} M_*^2, & \text{in the case of SMP.} \end{cases} \tag{59}$$

*Proof of Proposition 5* is completely similar to the one of Proposition 2 and is omitted.

In order to extract from (58) explicit efficiency estimates, we need to specify a step-size policy. In this respect, the advanced implementations offer more freedom than the basic ones. With the advanced implementation, at each iteration  $t$ , we ensure that  $r_t \leq 0$  and thus  $R_t \leq 0$ . This fact removes a technical complication from the analysis of the basic algorithm, namely we no longer need to ensure neither the martingale property of the random sums  $R_t$ , nor the monotonicity of the stepsizes. Therefore the step sizes in the advanced implementation can be far less restrictive than in the basic implementation. One option here is to use constant stepsize policy

$$\gamma_t = \sqrt{\frac{2\Omega}{N}} \cdot \begin{cases} \frac{1}{F_*}, & \text{in the case of SA} \\ \frac{1}{M_*}, & \text{in the case of SMP} \end{cases}, \quad 1 \leq t \leq N.$$

As it is easily seen, with this policy, (58) results in efficiency estimate [cf. (44)]

$$\forall t \geq 1, \quad \mathbf{E} \{ \epsilon_{\text{sad}}(z^t) \} \leq O(1) \begin{cases} F_* \sqrt{\frac{\Omega}{t}}, & \text{in the case of SA} \quad (a) \\ \mathcal{R}\mathcal{L} \sqrt{\frac{\Omega}{t}}, & \text{in the case of SMP} \quad (b) \end{cases} \quad (60)$$

Our preliminary experiments, however, suggest to equip the advanced implementations of SA and SMP with the online stepsize policy as follows. Let us set

$$\delta_t = \frac{\Theta^2}{t}, \quad S_t^* = \sum_{\tau=1}^t \delta_\tau \quad [ \leq \Theta^2(1 + \ln t) ] \quad (61)$$

and let us choose  $\gamma_\tau$  according to the “greedy” rule (the larger, the better) under the restriction that for all  $t = 1, 2, \dots$  it holds

$$R_t + S_t \leq S_t^*, \quad (*_t) \quad (62)$$

see (58). Specifically, assume that we have already carried out  $t - 1$  steps of the algorithm ensuring the relations  $(*_\tau)$ ,  $\tau \leq t - 1$ , and are about to define  $\gamma_t$  in order to carry out step  $t$  and to ensure  $(*_t)$ . When deciding on the value of  $\gamma_t$ , we already know the values of  $R_{t-1} \leq 0$  and  $S_{t-1}$ . Moreover we know in advance that whatever be our choice of  $\gamma_t > 0$ , we would have

$$R_t - R_{t-1} = r_t \leq 0, \quad S_t - S_{t-1} = s_t \leq \theta \gamma_t^2, \\ \theta = \begin{cases} \frac{F_*^2}{2}, & \text{in the case of SA} \\ \frac{M_*^2}{2} \leq 2\mathcal{L}^2\mathcal{R}^2, & \text{in the case of SMP} \end{cases}$$

[see (59)]. Thus, we can be sure that  $S_t + R_t \leq [S_{t-1} + R_{t-1}] + \theta \gamma_t^2$ , meaning that when choosing

$$\gamma_t = \sqrt{[S_t^* - S_{t-1} - R_{t-1}]/\theta} \quad (62)$$

we guarantee the validity of  $(*_t)$  and the inequality  $\gamma_t \geq \sqrt{\delta_t/\theta}$ . This observation combined with (58) and  $(*_N)$  implies that

$$\forall N \geq 1 : \quad \epsilon_{\text{sad}}(z^N) \leq \frac{\Theta^2/2 + R_N + S_N}{\sum_{\tau=1}^N \sqrt{\delta_\tau/\theta}} \leq \frac{O(1)\Theta^2(1 + \ln N)}{\sum_{\tau=1}^t \sqrt{\delta_\tau/\theta}} \\ \leq O(1)(1 + \ln N) \cdot \begin{cases} \Theta F_* N^{-1/2}, & \text{in the case of SA,} \\ \Theta \mathcal{R}\mathcal{L} N^{-1/2}, & \text{in the case of SMP.} \end{cases} \quad (63)$$

Observe that (63) is, within the logarithmic in  $N$  factor  $O(1)(1 + \ln N)$ , the same as the bound (60). In fact, we could somehow reduce this logarithmic gap by modifying  $S_t^*$ , but we do not think this is necessary; we may hope (and the experiments to be reported in Sect. 5 fully support this hope) that “in reality” the rule (62) is much better than it is stated by the above worst-case analysis. The rationale behind this hope is that while we indeed are conservative when thinking how large could  $S_t - S_{t-1}$  be, we account, to some extent, for the “past conservatism:” when  $S_{t-1} + R_{t-1}$  is essentially less than  $S_{t-1}^*$ ,  $\gamma_t$  as given by (62) is essentially larger than its lower bound used in the complexity analysis.

Finally, we remark that the major theoretical disadvantage of the efficiency estimate (63) as compared to (44) is much more serious than an extra log-factor. While with the basic implementation, in course of  $N$  steps the stochastic oracle is called  $O(1)N$  times, the number of oracle calls in course of  $N$  steps of the advanced implementation is random and can be much larger than  $O(1)N$ ; it is unclear why it should be  $O(1)N$  even on average. Though for the time being we cannot support the empirical evidence by a solid theoretical complexity analysis, in our experiments the advanced implementation by far outperformed its basic counterpart.

### 3.5 The favorable geometry case

We are about to present the “favorable geometry” case where we can point out the setup for SA/SMP which results in (nearly) *dimension-independent* efficiency estimates. Specifically, assume that

**[G.1]** The domain  $Z$  of  $(S)$  is a *subset* of the direct product  $Z^+ = B_1 \times \dots \times B_{p+q}$  of  $r = p + q$  “standard blocks” as follows:

- for  $1 \leq i \leq p$ ,  $B_i$  is the unit Euclidean ball in  $F_i = \mathbf{R}^{n_i}$ ;
- for  $1 \leq j \leq q$ ,  $B_{p+j}$  is a subset of the space  $F_{p+j}$  of  $n_{p+j} \times n_{p+j}$  ( $n_{p+j} > 1$ ) symmetric block-diagonal matrices of a given block-diagonal structure and is the *spectahedron* of  $F_{p+j}$ , that is, the set of all positive semidefinite matrices from  $F_{p+j}$  with unit trace.

In particular,  $B_{p+j}$  can be the standard simplex  $\{x \in \mathbf{R}_+^k : \sum_{\ell} x_{\ell} = 1\}$  (since the space of diagonal  $k \times k$  matrices can be naturally identified with  $\mathbf{R}^k$ ).

We equip  $F_i = \mathbf{R}^{n_i}$ ,  $i \leq p$ , with the standard Euclidean structure and the associated Euclidean norm  $\|\cdot\|_{(i)}$ , and  $F_{p+j}$ —with the Frobenius Euclidean structure and the trace-norm (the sum of singular values of a matrix)  $\|\cdot\|_{(p+j)}$ . In particular, the embedding space  $E = F_1 \times \dots \times F_r$  of  $Z^+$  becomes equipped with the direct product of the indicated Euclidean structures. Note that the norm  $\|\cdot\|_{(i,*)}$  conjugate to  $\|\cdot\|_{(i)}$  is either the norm  $\|\cdot\|_{(i)}$  itself (this is so when  $i \leq p$ ), or is the standard matrix norm (maximal singular value of a matrix) (this is so when  $i > p$ ). We denote a vector form on  $E$  as  $x = [x_1; \dots; x_r]$ , where  $x_{\ell}$  is the  $F_{\ell}$ -component of  $x$ .

**[G.2]** The decomposition  $Z = Z_1 \times Z_2 \subset E_1 \times E_2$  is compatible with the decomposition  $Z = B_1 \times \dots \times B_r$ , that is,  $E_1$  is the direct product of some of  $F_{\ell}$ ,  $1 \leq \ell \leq p + q$ , and  $E_2$  is the direct product of the remaining  $F_{\ell}$ . Besides this, we assume that  $Z$  intersects the relative interior of  $Z^+$ .



We refer to this case as to the one of *favorable geometry* and associate with this case the setup for SA and SMP as follows (cf. [13, Section 5]):

- The skew-symmetric linear mapping  $\mathcal{A}$  [see (3)] can be written down as

$$\mathcal{A}[x_1; \dots; x_r] = \left[ \sum_{j=1}^r A^{1j} x_j; \dots; \sum_{j=1}^r A^{rj} x_j \right],$$

where  $A^{ij}$  is a linear mapping from  $F_j$  to  $F_i$  and  $[A^{ij}]^* = -A^{ji}$ . We denote by  $L_{ij}$  an a priori upper bound on  $L_{ij}^* := \max_{x_j} \{ \|A^{ij} x_j\|_{(i,*)} : \|x_j\|_{(j)} \leq 1 \}$  such that  $L_{ij} = L_{ji}$ .<sup>5</sup> We assume that the symmetric matrix  $[L_{ij}]$  has no zero rows (this always can be enforced by replacing some of zero  $L_{ij}$ 's with small positive reals).

- Further, we set for  $1 \leq i \leq p$  and  $1 \leq j \leq q$ :

$$\omega_i(x_i) = \frac{1}{2} x_i^T x_i : B_i \rightarrow \mathbf{R}, \quad \Omega_i = \frac{1}{2},$$

$$\omega_{p+j}(x_{p+j}) = 2 \sum_{\ell=1}^{n_{p+j}} \lambda_\ell(x_{p+j}) \ln(\lambda_\ell(x_{p+j})) : B_{p+j} \rightarrow \mathbf{R}, \quad \Omega_{p+j} = 2 \ln(n_j),$$

where  $\lambda_\ell(u)$  are the eigenvalues of a symmetric matrix  $u$  taken with their multiplicities. It is known that  $\omega_\ell(\cdot)$  is a d.g.f. for  $B_\ell$  compatible with the norm  $\|\cdot\|_\ell$ ,  $1 \leq \ell \leq r$ .

- Finally, we define the norm  $\|\cdot\|$  on  $E$  and the d.g.f.  $\omega(\cdot)$  for  $Z$  according to

$$\mu_\ell = \frac{1}{\Omega_\ell} \frac{\sum_{j=1}^r L_{\ell j} \sqrt{\Omega_\ell \Omega_j}}{\sum_{i,j=1}^r L_{ij} \sqrt{\Omega_i \Omega_j}}, \quad \|[x_1; \dots; x_r]\| = \sqrt{\sum_{\ell=1}^r \mu_\ell \|x_\ell\|_{(\ell)}^2},$$

$$\omega(x) = \sum_{\ell=1}^r \mu_\ell \omega_\ell(x_\ell), \tag{64}$$

which results in

$$\Omega \leq 1, \quad \mathcal{R} \leq \Theta \leq \sqrt{2}, \quad \mathcal{L} = \sum_{i,j=1}^r L_{ij} \sqrt{\Omega_i \Omega_j}, \tag{65}$$

see [13, Section 5].

*Remark 1* From the results of [9] (see also Sect. A.2) it follows that the norm  $\|\xi\|_* = \sqrt{\sum_{\ell=1}^r \mu_\ell^{-1} \|\xi_\ell\|_{(i,*)}^2}$  is  $\varkappa$ -regular (see discussion in Sect. 3.3) with nearly dimension-independent  $\varkappa$ , namely,  $\varkappa = 7 \max_{1 \leq j \leq q} \ln(n_{p+j} + 1) + 1$ .

<sup>5</sup> The latter restriction is natural, since  $L_{ij}^* = L_{ji}^*$  due to  $[A^{ij}]^* = -A^{ji}$ .

Note that our motivating application of  $\ell_1$  minimization presented in Sect. 2.2.2 is of favorable geometry. The same is true for the low dimension approximation problem of Sect. 2.1.2; being a BSP rather than GBSP, this problem is well suited to illustrate the results we have obtained so far.

### 3.6 Illustration: low dimensional approximation via randomization

Passing in (9) from variable  $Q$  to variable  $R = d^{-1/2}Q$ , the problem reads

$$\min_{z_1 := \lambda \in Z_1 := \Delta_N} \max_{z_2 := R \in Z_2} \left[ d^{1/2} \text{Tr} \left( R \sum_{i=1}^N \lambda_i v_i v_i^T \right) \right], \tag{66}$$

$$Z_2 = \{R \in \mathbf{S}^n : 0 \preceq R \preceq d^{-1/2}I, \text{Tr}(R) = d^{1/2}\}.$$

We equip the embedding space  $E_1 = \mathbf{R}^N$  of  $Z_1$  with  $\|\cdot\|_1$ , and  $Z_1 = \Delta_N$ —with the entropy d.g.f. Further, we equip the embedding space  $E_2 = \mathbf{S}^n$  of  $Z_2$  with the Frobenius norm, and  $Z_2$  (which clearly is a subset of the unit ball of this norm)—with the Euclidean d.g.f.  $\frac{1}{2} \text{Tr}(z_2^2)$ . Taking into account that  $\|v_i\|_2 = 1$  for all  $i$ , it is immediately seen that we are in the Favorable Geometry case with  $\Omega_1 = 2 \ln(N)$ ,  $\Omega_2 = 1/2$ ,  $L_{12} = L_{21} = \sqrt{d}$  and  $L_{11} = L_{22} = 0$ .

Now assume that we want to solve (66) within a given accuracy  $\epsilon > 0$ . Consider  $t$ -step Basic implementation of SMP utilizing the distributions  $P_z^{(k)}$ ,  $k = \text{Ceil}(t \ln(N))$  (see Sect. 3.3.2) induced by the distributions  $P_z$  presented in Sect. 2.1.2, the stepsizes being given by (45). Taking into account Remark 1 and (65), we are in the situation of Corollary 1 with  $\Omega = O(1)$ ,  $\mathcal{L} = 2\sqrt{d} \ln(N)$ ,  $\sigma \leq O(1)\sqrt{\ln(N)/k} \mathcal{L} \leq O(1)\mathcal{L}/\sqrt{t}$ , so that (49) implies that  $\mathbf{E}\{\epsilon_{\text{sad}}(z^t)\} \leq O(1)\sqrt{d} \ln(N)/t$ . In particular, setting

$$t = t(\epsilon) = \text{Ceil} \left( O(1)\sqrt{d} \ln(N)/\epsilon \right) \tag{67}$$

with properly chosen  $O(1)$ , we ensure that  $\mathbf{E}\{\epsilon_{\text{sad}}(z^t)\} \leq \epsilon/2$ . Thus in course of running our algorithm, a solution of the required accuracy  $\epsilon$  will be built with probability  $\geq 1/2$ . Running our  $t(\epsilon)$ -step randomized procedure several times, until the first approximate solution with  $\epsilon_{\text{sad}}(z^t) \leq \epsilon$  is built (recall that  $\epsilon_{\text{sad}}(z^t)$  is observable on-line), we conclude that the probability *not* to find the desired approximate solution in  $mt(\epsilon)$  steps,  $m = 1, 2, \dots$ , is as small as  $2^{-m}$ .

Now let us look what, if any, is the gain of randomization. It is easily seen that in the case in question computing a value of the prox-mapping within machine precision costs  $O(n^3 + N)$  a.o. As a result, the best known so far complexity of solving (66) within accuracy  $\epsilon$  by any deterministic algorithm is, up to log-factors,  $\mathcal{C}_{\text{det}} = [n^3 + n^2N]\sqrt{d}\epsilon^{-1}$  a.o. According to Sect. 2.1.2, when sampling from  $P_z$ , after a “setup cost” of  $O(n^3 + dn^2 + N)$  a.o. is paid, generating a sample  $\zeta \sim P_z$  and computing  $F(\zeta)$  cost  $O(dnN)$  a.o. Thus, an iteration of the randomized method costs  $O(n^3 + dn^2 + dnN \underbrace{[t(\epsilon) \ln(N)]}_k)$  a.o., and the overall cost of an  $\epsilon$ -solution with this method,

again up to log-factors, is  $C_{\text{rand}} = \sqrt{d}\epsilon^{-1} \left[ n^3 + \sqrt{d}\epsilon^{-1}dnN \right]$  a.o. Assuming  $N \geq n$ , we get  $C_{\text{rand}}/C_{\text{det}} \leq O(1)[nN^{-1} + \epsilon^{-1}d^{3/2}n^{-1}]$ . For fixed  $\epsilon$ , this ratio tends to 0 as  $d, n, N$  grow in such a way that  $n/d^{3/2} \rightarrow +\infty$  and  $N/n \rightarrow +\infty$ .

### 4 Solving the generalized bilinear saddle point problem

Here we explain how a GBSP problem (11)–(12) can be reduced to a “small series” of BSP problems; the strategy to follow originates from [11]. From now on we assume, in addition to **A.1-2**, that we have an a priori upper bound  $\bar{\rho}$  on the optimal value  $\rho_*$  of (12). For example, it is immediately seen that when finding an  $\epsilon$ -solution to  $\ell_1$  minimization problem with  $\ell_p$  fit (Sect. 2.2.2) in the only nontrivial case  $\|b\|_p > \delta$  relation (15) implies that

$$\bar{\rho} := \|A\|_{1,p}[\|b\|_p - \delta]^{-1} \geq \rho_* := 1/\text{Opt}, \quad \|A\|_{1,p} = \max_j \|A_j\|_p, \quad (68)$$

where  $A_1, \dots, A_n$  are the columns of  $A$ .

For the sake of definiteness, we assume that we are in the Favorable Geometry case, and that the decomposition  $Z = Z_{11} \times Z_{12} \times Z_2 \subset E$ , see (13), is compatible with the decomposition  $E = F_1 \times \dots \times F_r$ , that is, the embedding spaces of  $Z_{11}, Z_{12}$  and  $Z_2$  are products of some of  $F_\ell$ 's. To save space, we restrict ourselves with the SMP algorithm; modifications in the case of SA are straightforward.

*The algorithm* solves the problem of interest (12) by applying to  $\text{SV}(\cdot)$  a Newton-type root finding routine, with (approximate) first order information on  $\text{SV}$  at a point  $\rho$  given by SMP as applied to the BSP problem specifying  $\text{SV}(\rho)$ . Specifically, the algorithm works stage by stage. At a stage  $s$ , we have at our disposal an upper bound  $\rho_s$  on  $\rho_*$  and a piecewise linear function  $\ell_{s-1}(\rho)$  which underestimates  $\text{SV}(\cdot)$ :

$$\text{SV}(\rho) \geq \ell_{s-1}(\rho) \quad \forall \rho \geq 0.$$

here  $\rho_1 = \bar{\rho}, \ell_0 \equiv -\infty$ . At a stage, we apply SMP to the BSP problem

$$\text{SV}(\rho_s) = \min_{z_1 \in Z_1} \max_{z_2 \in Z_2} \phi^{\rho_s}(z_1, z_2) \quad (\mathcal{S}_s)$$

namely, act as follows.

A. We start stage  $s$  with building the setup for SMP as explained in Sect. 3.5. The affine operator associated with  $(\mathcal{S}_s)$  is

$$\begin{aligned} F^{\rho_s}(z_1 = [z_{11}; z_{12}], z_2) &= \Phi(z_1, z_2) + \rho_s \Psi(z_1, z_2) \\ &= \left[ \begin{array}{l} [a_{11} + B^*z_2; \rho_s(a_{12} + C^*z_2)] \\ -b - Bz_{11} - \rho_s(c + Cz_{12}) \end{array} \right], \end{aligned}$$

see (11), (13). In matrix  $\mathcal{A} = \mathcal{A}_s$  of the linear part of  $F^{\rho_s}$ , some blocks  $A^{ij}$  are independent of  $\rho_s$ , while the remaining blocks are proportional to  $\rho_s$ . Consequently, the

Lipschitz constant of  $F^{\rho_s}$  as given by (65) is

$$\mathcal{L} = \mathcal{L}(\rho_s) = \mathcal{M} + \rho_s \mathcal{N}, \quad \mathcal{M}, \mathcal{N} \geq 0. \tag{69}$$

Observe that by Remark 1, the regularity parameter of the norm  $\|\cdot\|_* = \|\cdot\|_*^{(s)}$  conjugate to the norm  $\|\cdot\| = \|\cdot\|^{(s)}$  participating in the SMP setup for  $s$ th stage does not exceed

$$\varkappa = 7 \ln(N + 1) + 1, \tag{70}$$

where  $N$  is the largest size of the spectahedron blocks, if any (otherwise  $N = 0$ ), participating in  $Z$ .

*B.* We apply to  $(\mathcal{S}_s)$  either the basic, or the advanced implementation of the SMP. When running the basic SMP, we use the distributions  $P_z^{(k)}$ , see Sect. 3.3 (here  $k \geq 1$  is a parameter of the construction) and the online stepsize policy (45), where we set  $\mathcal{L} = \mathcal{L}_s := \mathcal{M} + \rho_s \mathcal{N}$  and  $\Omega = 1$  [see (65)]. In addition, we restart the basic SMP every

$$N_s(\rho_s \epsilon) := \text{Ceil} \left( \max \left[ \frac{102 \mathcal{L}_s}{\rho_s \epsilon}, \left( \frac{272 \mathcal{L}_s}{\rho_s \epsilon} \right)^2 \frac{\varkappa}{k} \right] \right) \tag{71}$$

steps, see below; here  $\varkappa$  is given by (70). When  $(\mathcal{S}_s)$  is solved by the advanced SMP, we use the online stepsize policy (61)–(62), with  $\Theta = \sqrt{2}$  in (61).

*B.I.* Let  $z^{ts} = [z_1^{ts}, z_2^{ts}]$  be the approximate solution to  $(\mathcal{S}_s)$  generated after  $t$  steps of stage  $s$ ; recall that along with this solution, we have at our disposal the quantities

$$\begin{aligned} \bar{\phi}^{ts} &= \max_{z_2 \in Z_2} \phi^{\rho_s}(z_1^{ts}, z_2) = \nu + \langle a_{11}, z_{11}^{ts} \rangle + \rho_s [\chi + \langle a_{12}, z_{12}^{ts} \rangle] \\ &\quad + \min_{z_2 \in Z_2} \langle z_2, b + \rho_s c + B z_{11}^{ts} + \rho C z_{12}^{ts} \rangle, \\ \underline{\phi}^{ts} &= \min_{z_1 \in Z_1} \phi^{\rho_s}(z_1, z_2^{ts}) = \overbrace{\nu + \langle b, z_2^{ts} \rangle}^{p_{ts}} + \min_{z_{11} \in Z_{11}} \langle a_{11} + B^* z_2^{ts}, z_{11} \rangle \\ &\quad + \rho_s \left[ \overbrace{\chi + \langle c, z_2^{ts} \rangle}^{q_{ts}} + \min_{z_{12} \in Z_{12}} \langle a_{12} + C^* z_2^{ts}, z_{12} \rangle \right] \end{aligned} \tag{72}$$

[cf. (29) and see (11), (13)]. We set

$$u^{ts} = \min_{\tau \leq t} \bar{\phi}^{\tau s}, \quad \ell^{ts} = \max_{\tau \leq t} \underline{\phi}^{\tau s}, \quad \ell_{ts}(\rho) = \max \left[ \ell_{s-1}(\rho), \max_{1 \leq \tau \leq t} [p_{\tau s} + q_{\tau s} \rho] \right].$$

Note that  $u^{ts}$  is a nonincreasing in  $t$  upper bound on  $\text{SV}(\rho_s)$ ,  $\ell^{ts}$  is a nondecreasing in  $t$  lower bound on  $\text{SV}(\rho_s)$ , and  $\ell_{ts}(\rho)$  underestimates  $\text{SV}(\rho)$  for all  $\rho \geq 0$ . In addition,

$\ell_{ts}(\rho_s) \geq \ell^{ts}$ . Note also that after  $t$  steps we have at our disposal vectors  $w_1^{ts} \in Z_1$ ,  $w_2^{ts} \in Z_2$  such that

$$\max_{z_2 \in Z_2} \phi^{\rho_s}(w_1^{ts}, z_2) = u^{ts} \leq \bar{\phi}^{ts}, \quad \min_{z_1 \in Z_1} \phi^{\rho_s}(z_1, w_2^{ts}) = \ell^{ts} \geq \underline{\phi}^{ts},$$

meaning that  $w^{ts} = [w_1^{ts}; w_2^{ts}]$  is a feasible solution to  $(S_s)$  and  $\epsilon_{\text{sad}}(w^{ts}) = u^{ts} - \ell^{ts} \leq \bar{\phi}^{ts} - \underline{\phi}^{ts} = \epsilon_{\text{sad}}(z^{ts})$ .

B.2. We proceed with solving  $(S_s)$  until one of the following three situations occurs:

- (A) We get  $u^{ts} \leq \epsilon \rho_s$ . In this case we terminate with the claim that  $\rho_s, w_1^{ts}$  is the desired  $\epsilon$ -solution to (11)–(12).
- (B) We get  $\ell^{ts} \geq \frac{3}{4}u^{ts}$ . When it happens, we set

$$\rho_{s+1} = \max \{ \rho : \ell_{ts}(\rho) \leq 0 \}, \quad \ell_s(\cdot) \equiv \ell_{ts}(\cdot) \tag{73}$$

and pass to the stage  $s + 1$ .

- (C) The iteration count  $t$  becomes a multiple of  $N_s(\epsilon)$ . When it happens *and if the basic implementation of SMP is used*, we restart SMP and proceed to step  $t + 1$  of stage  $s$  (that is, the next iterate of stage  $s$  will be  $z_\omega$ , the subsequent approximate solutions will be weighted sums of the points  $w$  generated after the restart, etc.) If the advanced implementation of SMP is used, we do not restart the algorithm and proceed as at all other steps.

**Theorem 1** *When solving a Generalized Bilinear Saddle Point problem (11)–(12) within the accuracy  $\epsilon > 0$  by the outlined algorithm:*

- (i) *The algorithm terminates in finite time with probability 1, and the resulting solution is an  $\epsilon$ -solution, as defined in Sect. 2.2, to the GBSP problem in question;*
- (ii) *The number of stages does not exceed the quantity  $O(1) \ln \left( \frac{\|\phi\|_\infty + \bar{\rho} \|\psi\|_\infty}{\epsilon \rho_*} + 2 \right)$ , where  $\|\phi\|_\infty = \max_{z \in Z} |\phi(z)|$ ,  $\|\psi\|_\infty = \max_{z \in Z} |\psi(z)|$ ,  $\rho_*$  is the optimal value in the problem (11)–(12), and  $\bar{\rho}$  is an a priori upper bound on  $\rho_*$ , see the beginning of Sect. 4.*
- (iii) *The (random) number  $N_s$  of steps at every stage  $s$  of the basic implementation satisfies the relation*

$$\text{Prob}\{N_s \geq m N_s(\rho_s \epsilon)\} \leq 2^{-m}, \quad m = 1, 2, \dots \tag{74}$$

with  $N_s(\rho_s \epsilon)$  given by (71). Besides this,

$$N_s(\rho_s \epsilon) \leq O(1) \frac{\mathcal{L} + \rho_* \mathcal{N}}{\rho_* \epsilon} \left[ 1 + \chi \frac{\mathcal{L} + \rho_* \mathcal{N}}{\rho_* \epsilon k} \right] \tag{75}$$

with  $\chi$  given by (70).

The number of steps at every stage of the advanced implementation of the algorithm does not exceed

$$N_{\text{adv}}(\epsilon) = O(1) \left[ \frac{\mathcal{M} + \rho_* \mathcal{N} + 2\epsilon \rho_*}{\epsilon \rho_*} \ln \left( \frac{\mathcal{M} + \rho_* \mathcal{N} + 2\epsilon \rho_*}{\epsilon \rho_*} \right) \right]^2. \tag{76}$$

*Proof 1<sup>0</sup>.* From the description of the method it follows that

$$\forall t, s \geq 1, \forall \rho \geq 0 : u^{ts} \geq \text{SV}(\rho_s) \geq \ell^{ts}, \quad \ell_{ts}(\rho) \leq \text{SV}(\rho), \quad \ell^{ts} \leq \ell_{ts}(\rho_s). \tag{77}$$

Let us prove by induction in  $s$  that  $\rho_* \leq \rho_s \leq \rho_1$ . The base  $s = 1$  is evident. Now let  $\rho_* \leq \rho_s \leq \rho_1$ , and let stage  $s + 1$  take place. When passing from stage  $s$  to stage  $s + 1$ , we are in the case (B) and thus have  $u^{ts} > \epsilon \rho_s$ ,  $\ell^{ts} \geq \frac{3}{4}u^{ts} > \frac{3}{4}\epsilon \rho_s$ , whence, in view of (77),

$$\ell_s(\rho_s) = \ell_{ts}(\rho_s) \geq \ell^{ts} \geq \frac{3}{4} \max[\epsilon \rho_s, \text{SV}(\rho_s)] \quad \text{thus} \quad \ell_s(\rho_s) > 0. \tag{78}$$

This combines with  $\ell_{ts}(\rho_*) \leq \text{SV}(\rho_*) \leq 0$  and convexity of  $\ell_{ts}(\cdot)$  to imply that  $\rho_* \leq \rho_{s+1} < \rho_s$ . Induction is complete.

Since  $\rho_s \geq \rho_*$ ,  $u^{ts}$  is an upper bound on  $\text{SV}(\rho_s)$  and  $u^{ts} \geq \bar{\phi}^{\rho_s}(w_1^{ts})$ , we conclude that if the algorithm terminates at stage  $s$ , then the result  $\rho_s, w_1^{ts}$  is an  $\epsilon$ -solution to the GBSP in question.

*2<sup>0</sup>.* Let us prove (ii). The reasoning to follow goes back to [11]; we reproduce it here to make the paper self-contained. Let  $s$  be such that the stage  $s + 1$  takes place, and let  $u_s$  be the last bound  $u^{ts}$  built at stage  $s$ . Observe that

$$\frac{3}{4}\epsilon \rho_s < \frac{3}{4}u_s \leq \ell_s(\rho_s) \leq \text{SV}(\rho_s) \leq u_s. \tag{79}$$

Since the convex function  $\ell_s(\rho)$  is nonpositive at  $\rho = \rho_{s+1}$  and is  $\geq \frac{3}{4}u_s > 0$  at  $\rho = \rho_s > \rho_{s+1}$ , we have  $g_s := \ell'_s(\rho_s) > 0$  and

$$\rho_s - \rho_{s+1} \geq \ell_s(\rho_s)/g_s \geq \frac{3}{4}u_s/g_s. \tag{80}$$

Now assume that  $s > 1$  is an intermediate step, i.e., it is such that the stage  $s + 1$  also takes place. Applying (80) and (79) to  $s - 1$  in the role of  $s$ , we get  $\rho_{s-1} - \rho_s \geq \frac{3}{4}u_{s-1}/g_{s-1}$  and  $\frac{3}{4}u_s \leq \ell_s(\rho_s)$ , whence, by convexity of  $\ell_s(\cdot)$  and in view of (77), we have

$$u_{s-1} \geq \text{SV}(\rho_{s-1}) \geq \ell_s(\rho_{s-1}) \geq \ell_s(\rho_s) + g_s(\rho_{s-1} - \rho_s) \geq \frac{3}{4}u_s + g_s \frac{3}{4} \frac{u_{s-1}}{g_{s-1}}.$$

Consequently,  $\frac{4}{3}u_{s-1} \geq u_s + \frac{g_s u_{s-1}}{g_{s-1}}$ , or equivalently  $\frac{u_s}{u_{s-1}} + \frac{g_s}{g_{s-1}} \leq \frac{4}{3}$ , whence  $\frac{u_s g_s}{u_{s-1} g_{s-1}} \leq (1/4)(4/3)^2 = 4/9$ . It follows that

$$\sqrt{u_s g_s} \leq (2/3)^{s-1} \sqrt{u_1 g_1}. \tag{81}$$

We have  $\ell_s(\rho_*) \leq \text{SV}(\rho_*) = 0$ ,  $\ell_s(\rho_s) \geq \frac{3}{4}u_s$  [see (79)] and  $\ell_s(\rho_s) - \ell_s(\rho_*) \leq g_s(\rho_s - \rho_*)$  (by convexity of  $\ell_s(\cdot)$ ), whence  $g_s \geq \frac{3}{4}u_s(\rho_s - \rho_*)^{-1} \geq \frac{3}{4\rho_1}u_s$ , and (81) implies that

$$u_s \leq (2/3)^{s-1} \sqrt{u_1 g_1} \sqrt{4\rho_1/3}. \tag{82}$$

Now,  $g_1 = \ell'_1(\rho_1)$  and  $\ell_1(\rho) \leq \text{SV}(\rho) \leq \|\phi\|_\infty + \rho\|\psi\|_\infty$ , and  $g_1 \leq \|\psi\|_\infty$ , and clearly  $u_1 \leq \|\phi\|_\infty + \rho_1\|\psi\|_\infty$ . At the same time,  $u_s > \epsilon\rho_s \geq \epsilon\rho_*$ , so that (82) implies that  $\epsilon\rho_* \leq (2/3)^{s-1}[\|\phi\|_\infty + \rho_1\|\psi\|_\infty]$ . The resulting upper bound on  $s$  implies (ii).

**3<sup>0</sup>.** Let us prove (iii). Assume, first, that Basic SMP with stepsizes (45) is used. From the description of the algorithm it follows that at every stage  $s$ , before termination of the stage, the residual of current approximate solutions  $w^{ts}$  is  $> \frac{1}{4}\epsilon\rho_s$  (since  $u^{ts} > \epsilon\rho_s$  and  $\ell^{ts} < \frac{3}{4}u^{ts}$ ). It follows that in order to prove (74), it suffices to verify that when applying to  $(\mathcal{S}_s)$   $N = N_s(\rho_s\epsilon)$ -step Basic SMP, we have  $\epsilon_{\text{sad}}(z^N) \leq \epsilon\rho_s/4$  with probability  $\geq 1/2$ ; to this end, it is enough to verify that the expectation of  $\epsilon_{\text{sad}}(z^N)$  is  $\leq \epsilon\rho_s/8$ . By Corollary 1, this expectation is  $\leq \alpha := 7\Omega\mathcal{L}_s/N + 6\sqrt{\Omega}\sigma/\sqrt{N}$ , where  $\sigma^2 = \sup_{z \in Z} \mathbf{E}_{\zeta \sim P_z^{(k)}} \{\|\mathcal{A}(\zeta - z)\|_*^2\} \leq \frac{\kappa}{k}(2\mathcal{L}_s\Theta)^2$  [see (36) and discussion in Sect. 3.3.2], that is,  $\sigma \leq 2\Theta\mathcal{L}_s\sqrt{\kappa/k}$ . This inequality combines with the relations  $\Omega = 1$ ,  $\Theta = \sqrt{2\Omega}$  and the definition (71) of  $N = N_s(\rho_s\epsilon)$  to imply the desired bound  $\alpha \leq \epsilon\rho_s/8$ . We have proved (74); (75) is readily given by (71) and the relation  $\rho_s \geq \rho_*$ .

For the advanced implementation of SMP, similar reasoning based on the bound (63) with  $\mathcal{L} = \mathcal{M} + \rho_s\mathcal{N}$  justifies (76).

**4<sup>0</sup>.** Combining (ii), (iii) and the concluding claim in item 1<sup>0</sup> above, we arrive at (i). □

*The case of  $\ell_1$  minimization.* In the case of  $\ell_1$  minimization problems with uniform and  $\ell_2$  fits, Theorem 1 as applied to the basic implementation of SMP with  $k = 1$ , initialized according to (68), after completely straightforward computations implies the complexity bounds stated in Proposition 1. The preprocessing mentioned in item (ii) of Proposition 1 is as follows: we choose an  $m \times m$  orthogonal matrix  $U$  with moduli of entries not exceeding  $O(1)/\sqrt{m}$  and such that multiplication of a vector by  $U$  takes  $O(m \ln m)$  operations (e.g.,  $U$  can be the matrix of the Cosine Transform). We then draw at random a  $\pm 1$  vector  $\xi$  from the uniform distribution on the vertices of the unit  $m$ -dimensional box and pass from the data  $[A, b]$  to the data

$$[A' = U\text{Diag}\{\xi\}A, \quad b' = U\text{Diag}\{\xi\}b],$$

thus obtaining an equivalent reformulation of the problem of interest. Note that this preprocessing costs  $O(1)mn \ln(m)$  operations. We clearly have  $\|A'\|_{1,2} = \|A\|_{1,2}$ . Applying Hoeffding’s inequality (see [8]), it is immediately seen that for every tolerance  $\chi \in (0, 1/2)$  with probability  $\geq 1 - \chi$  one has  $\|A'\|_{1,\infty} < O(1)\sqrt{\ln(mn/\chi)}m^{-1/2}\|A\|_{1,2}$ , that is,  $\Gamma(A') \leq O(1)\sqrt{\ln(mn/\chi)}$ , as stated in Proposition 1.

### 5 Numerical results

Below we report on a series of numerical experiments aimed at comparing the performances of the Stochastic Mirror Prox algorithm SMP (in its advanced implementation) and its prototype—deterministic mirror prox algorithm (DMP) proposed in [13].<sup>6</sup> The algorithms were tested on the GBSP problems of  $\ell_1$  minimization with uniform and  $\ell_2$  fits, see Sect. 2.2.2.

*Test problems* we use originate from compressive sensing. Specifically, given the sizes  $m, n$  of a test problem, we picked at random an  $m \times n$  matrix  $B$  with i.i.d. entries taking values  $\pm 1$  with probabilities 0.5, and a sparse (with  $\text{Ceil}(\sqrt{m})$  nonzero entries randomly generated from standard Normal distribution) “true signal”  $x_*$  normalized to have  $\|x_*\|_1 = 1$ , thus giving rise to the test problem

$$\text{Opt}_p = \min_x \{ \|x\|_1 : \|Ax - y\|_p \leq \delta \}, \quad A = m^{-1/p}B, \quad y = Ax_* + \xi \quad (P_p)$$

where  $p = \infty$  (uniform fit) or  $p = 2$  ( $\ell_2$  fit). The “observation noise”  $\xi$  was chosen at random (each entry is from an i.i.d. standard Normal distribution) and then normalized to have  $\|\xi\|_p = \delta$ , thus making sure that the true solution  $x_*$  is feasible to  $(P_p)$ . Our goal is to solve  $(P_p)$  within accuracy  $\epsilon$ , i.e., to find  $x_\epsilon$  satisfying  $\|x_\epsilon\|_1 \leq \text{Opt}_p$  and  $\|Ax_\epsilon - y\|_p \leq \delta + \epsilon$ . In all our experiments,  $\delta = 0.005$  and  $\epsilon = 0.0025$  were used.

*Implementation of the algorithms* The GBSP reformulations of problems  $(P_p)$  were solved by SMP (in advanced implementation) and DMP according to the scheme presented in Sect. 4. In the case  $p = \infty$  of uniform fit, both SMP and DMP used the GBSP problem reformulation given by (19). In the case  $p = 2$  of  $\ell_2$  fit, SMP used the GBSP reformulation (18), while DMP was applied to the GBSP problem stemming directly from (16) with  $p = 2$ , namely, given by

$$\phi^p(z_1, z_2) = z_2^T (AJ_n z_1 - \rho b) - \rho \delta, \quad Z_1 = Z_{11} = \Delta_{2n}, \quad Z_2 = \{\|z_2\|_2 \leq 1\}. \quad (83)$$

The rationale here is that the GBSP given by (83) “by itself” is easier than the GBSP given by (18): an  $\epsilon$ -solution to the latter problem induces straightforwardly an  $\epsilon$ -solution to the former one, but not vice versa. As a compensation, the problem (18), in

<sup>6</sup> DMP is nothing but SMP with precise information (i.e.,  $P_z$  is the unit mass sitting at  $z$ ) and on-line stepsize policy described in [13, Section 6].



contrast to (83), is better suited for randomization.<sup>7</sup> The latter fact, which is crucial for SMP, is irrelevant for DMP, this is why we apply this algorithm to the GBSP given by (83). In order to make a fair comparison, when running SMP for  $\ell_2$ -fit, we terminate the run based on the  $\ell_2$ -residual of the solution.

In our implementations, we have tested different policies for choosing the starting point at each stage and different choices of the distance generating function (d.g.f.) for the simplexes. Specifically, along with the entropy d.g.f. discussed in Sect. 3.5, we tested the power d.g.f.  $\omega(x) = \frac{e}{\kappa(1+\kappa)} \sum_{i=1}^n x_i^{1+\kappa} : \{x \in \mathbf{R}_+^n : \sum_i x_i \leq 1\} \rightarrow \mathbf{R}$ , with  $\kappa = \frac{1}{\ln(n)}$ ; the theoretical complexity bounds associated with this choice of d.g.f. coincide, within absolute constant factors, with those for the entropy. The best policies we ended up with are as follows:

- for SMP: entropy d.g.f., restarts from the  $\omega$ -center of  $Z$  (“C00E” implementation);
- for DMP, in the case of uniform fit: power d.g.f., restarts from the convex combination of the best (with the smallest  $\epsilon_{\text{sad}}$ ) point found so far and the  $\omega$ -center of  $Z$ , the weights being 0.25 and 0.75, respectively (“B25P” implementation);
- for DMP, in the case of  $\ell_2$ -fit: power d.g.f., restarts from the convex combination of the last search point of the previous stage and the  $\omega$ -center of  $Z$ , the weights being 0.75 and 0.25, respectively (“L75P” implementation).

When implementing SMP, we utilized the option, discussed in Sect. 3.3, of building an estimate  $F(\zeta)$  of  $F(z)$  by generating  $k$  samples  $\zeta^\ell \sim P_z, \ell = 1, \dots, k$ , and setting  $\zeta = \frac{1}{k} \sum_{\ell=1}^k \zeta^\ell$ . The “multiplicity”  $k$  was set to 40 for small instances and 100 for large (those with at least  $10^8$  nonzeros in  $A$ ) instances.

The MATLAB 7.10.0 implementation of the algorithms was executed on an eight-core machine with two quad-core Intel Xeon E5345 CPU@2.33 GHz, 8 MB L2 cache per quad-core chip and 12 GB FB-DIMM total RAM (the computations were running single-core and single-threaded).

*The results, I* In order to avoid too time-consuming experimentation, we primarily dealt with “moderate size” test problems. These problems were split into four groups according to the total number of nonzeros in  $A$  ( $2 \cdot 10^6, 8 \cdot 10^6, 32 \cdot 10^6, 128 \cdot 10^6$ ). Every group was further split into two subgroups according to the ratio  $n : m$  (8 and 2). For every one of the resulting pairs  $(m, n)$ , we generated 5 instances of problem  $(P_2)$  and 5 instances of problem  $(P_\infty)$  and solved them by DMP and SMP. Thus, the methods were compared on totally 80 problem instances split into 16 series of 5 experiments each, with common for all experiments of a series sizes  $m, n$  and the value of  $p$ . The results are presented in Tables 1 (uniform fit) and 2 ( $\ell_2$  fit). For every series of 5 experiments, we present the corresponding minimal, maximal and average values of several performance characteristics, specifically

- CPU—the CPU time (measured in seconds (s)) of the entire computation

<sup>7</sup> Indeed, in the second problem all nontrivial matrix-vector multiplications required to compute  $F^\rho(z)$  are multiplications of vectors from the  $\ell_1$ -balls by  $A$  and  $A^T$ ; since a vector from  $\ell_1$ -ball is the expectation of an extremely sparse (just one nonzero entry) random vector taking values in the same ball, the required matrix-vector multiplications admit cheap randomized versions. In the first problem, some of the required matrix-vector multiplications involve vectors from the  $\|\cdot\|_2$ -ball, and such a vector typically cannot be represented as the expectation of a sparse random vector taking values in the ball.

**Table 1** Numerical results for  $\ell_1$ -minimization with  $\|\cdot\|_\infty$ -fit

Sizes	DMP		SMP			I <sup>a</sup>	II <sup>b</sup>
	Calls	CPU (s)	Calls	FCalls	CPU (s)		
<b>500 × 4,000</b>							
Mean (C00E)	2,661.6	106.6	10,511.0	236.5	57.2	11.89	1.98
Min(C00E)	1,683.0	50.0	8,159.0	183.6	34.0	6.91	1.16
Max (C00E)	4,395.0	179.4	11,783.0	265.1	83.4	23.94	4.14
Mean (B25P)	1,453.4	104.1				6.15	1.89
<b>1,000 × 2,000</b>							
Mean (C00E)	1,830.8	64.0	10,568.8	158.5	42.9	11.69	1.54
Min (C00E)	1,344.0	41.0	8,434.0	126.5	28.8	7.82	1.02
Max (C00E)	2,507.0	91.5	11,576.0	173.6	70.4	15.83	2.02
Mean (B25P)	1,530.6	97.9				9.64	2.48
<b>1,000 × 8,000</b>							
Mean (C00E)	2,338.0	227.9	12,406.6	139.6	113.2	16.68	1.99
Min (C00E)	1,453.0	119.4	11,579.0	130.3	88.2	11.15	1.27
Max (C00E)	2,739.0	370.2	13,895.0	156.3	168.9	18.99	2.39
Mean (B25P)	1,545.6	248.9				11.08	2.30
<b>2,000 × 8,000</b>							
Mean (C00E)	2,691.6	227.6	12,922.8	96.9	74.5	27.93	3.10
Min (C00E)	1,132.0	97.7	10,934.0	82.0	56.6	12.24	1.37
Max (C00E)	3,355.0	313.1	15,632.0	117.2	88.8	35.46	4.25
Mean (B25P)	1,426.4	207.8				14.74	2.84
<b>2,000 × 16,000</b>							
Mean (C00E)	2,384.6	494.2	13,174.8	74.1	184.9	32.30	2.68
Min (C00E)	2,288.0	486.3	11,735.0	66.0	174.4	29.78	2.53
Max (C00E)	2,491.0	505.5	14,729.0	82.9	195.3	34.66	2.84
Mean (B25P)	1,575.2	533.7				21.41	2.89
<b>4,000 × 8,000</b>							
Mean (C00E)	2,923.6	798.7	19,750.2	74.1	228.4	39.42	3.30
Min (C00E)	2,032.0	407.6	17,262.0	64.7	159.0	28.86	2.34
(C00E)	3,895.0	1,539.7	22,945.0	86.0	343.1	48.61	4.49
Mean (B25P)	1,554.6	576.2				21.12	2.63
<b>4,000 × 32,000</b>							
Mean (C00E)	2,482.8	2,054.3	11,973.2	84.2	515.8	29.47	3.98
Min (C00E)	1,826.0	1,448.9	11,331.0	79.7	499.9	22.39	2.90
Max (C00E)	3,479.0	2,904.2	12,715.0	89.4	525.0	42.65	5.70
Mean (B25P)	1,604.8	1,736.3				19.19	3.36
<b>8,000 × 16,000</b>							
Mean (C00E)	2,680.4	2,227.7	12,474.6	58.5	375.0	45.78	5.92
Min (C00E)	2,297.0	1,890.1	11,493.0	53.9	341.9	41.12	5.44
Max (C00E)	3,177.0	2,609.0	13,759.0	64.5	408.8	49.26	6.48
Mean (B25P)	1,615.8	1,752.7				27.57	4.63

<sup>a</sup> Calls, DMP    <sup>b</sup> CPU, DMP  
 FCalls, SMP    CPU, SMP

**Table 2** Numerical results for  $\ell_1$ -minimization with  $\|\cdot\|_2$ -fit

Sizes	DMP		SMP			I <sup>a</sup>	II <sup>b</sup>
	Calls	CPU (s)	Calls	FCalls	CPU (s)		
<b>500 × 4,000</b>							
Mean (C00E)	579.8	21.0	4,771.6	106.7	24.6	5.91	0.93
Min (C00E)	410.0	14.5	3,412.0	76.3	16.9	3.18	0.49
Max (C00E)	722.0	40.3	6,868.0	153.5	36.0	8.40	1.94
Mean (L75P)	287.8	16.1				2.95	0.70
<b>1,000 × 2,000</b>							
Mean (C00E)	553.0	19.0	3,910.8	54.8	13.6	10.73	1.47
Min (C00E)	463.0	9.1	3,315.0	46.4	11.5	5.68	0.52
Max (C00E)	664.0	30.1	5,890.0	82.5	17.4	13.56	2.34
Mean (L75P)	282.4	14.1				5.44	1.07
<b>1,000 × 8,000</b>							
Mean (C00E)	617.0	56.6	5,148.8	57.5	50.7	11.25	1.17
Min (C00E)	486.0	34.7	3,745.0	41.9	36.1	7.68	0.74
Max (C00E)	794.0	87.1	6,050.0	67.6	64.8	18.35	1.93
Mean (L75P)	318.8	40.9				5.84	0.86
<b>2,000 × 8,000</b>							
Mean (C00E)	634.8	39.8	5,853.6	41.0	47.2	15.94	0.86
Min (C00E)	487.0	30.0	3,926.0	27.5	33.1	11.17	0.59
Max (C00E)	796.0	51.0	6,869.0	48.1	54.0	20.49	1.12
Mean (L75P)	318.8	25.9				8.05	0.58
<b>2,000 × 16,000</b>							
Mean (C00E)	531.8	150.7	5,055.6	28.3	90.0	19.88	1.80
Min (C00E)	438.0	108.3	3,947.0	22.1	60.2	11.64	0.87
Max (C00E)	608.0	180.3	6,736.0	37.6	125.1	24.80	2.49
Mean (L75P)	346.0	110.6				12.74	1.28
<b>4,000 × 8,000</b>							
Mean (C00E)	675.2	138.5	6,504.6	22.8	101.7	29.71	1.36
Min (C00E)	531.0	99.1	5,868.0	20.5	83.3	22.71	0.99
Max (C00E)	810.0	193.6	7,143.0	25.0	113.9	34.52	1.70
Mean (L75P)	346.4	86.3				15.21	0.85
<b>4,000 × 32,000</b>							
Mean (C00E)	672.2	486.0	5,613.4	39.2	287.2	17.66	1.74
Min (C00E)	506.0	382.5	3,418.0	23.9	197.2	12.08	1.26
Max (C00E)	817.0	579.1	6,611.0	46.2	336.4	22.57	2.15
Mean (L75P)	355.4	311.6				9.39	1.12
<b>8,000 × 16,000</b>							
Mean (C00E)	592.4	591.4	5,815.0	25.4	177.6	24.15	3.51
Min (C00E)	509.0	472.4	3,765.0	16.5	117.3	16.56	2.36
Max (C00E)	696.0	798.1	7,038.0	30.8	214.1	30.90	5.06
Mean (L75P)	329.8	360.2				13.38	2.10

<sup>a</sup> Calls, DMP   <sup>b</sup> CPU, DMP  
 FCalls, SMP   CPU, SMP

- **Calls**—the total number of computations of the values of  $F$
- **FCalls**—the equivalent number of calls to the deterministic oracle for the randomized algorithm. This quantity is defined as follows. For DMP, computing a value of  $F$  at a point reduces to a pair of matrix-vector multiplications, one involving  $A$  and the other one involving  $A^T$ ; the cost of this computation is  $2mn$  operations. For SMP invoked with multiplicity  $k$  (see above), the computation of (an unbiased estimate of)  $F(z)$  requires multiplying one vector with  $\leq k$  nonzero entries by  $A$ , and another vector with  $\leq k$  nonzero entries by  $A^T$ , the total cost of these two computations being  $k(m+n)$  operations. Thus, the “deterministic equivalent” of the randomized computation of  $F$  used by SMP is  $\frac{k(m+n)}{2mn}$ . The quantity **FCalls** represents the equivalent number of calls to the deterministic oracle that we could afford for the same total computational cost involved with the queries to the stochastic oracle needed to solve the problem by SMP.

The data in Tables 1 and 2 suggest the following interpretations:

1. As the sizes of instances grow, the randomized algorithm eventually outperforms its deterministic counterpart in terms of the CPU time, and the corresponding “savings” grow with the size  $m \times n$  of the instance, and for instances of a given size—grow as the ratio  $n/m$  decreases. Both phenomena are quite natural: the larger is  $mn$  and the smaller is  $n/m \geq 1$  for a given  $mn$ , the smaller is the deterministic equivalent  $k \frac{m+n}{2mn}$  of a randomized computation of  $F$ .
2. Even for our “not too large” test problems, the savings stemming from randomization can be quite significant: for the  $8,000 \times 16,000$  instances, SMP is, at average, nearly 4.6 times faster than the best version of DMP for problems with uniform fit and 2.1 times faster than DMP for problems with  $\ell_2$  fit.

When interpreting the CPU time data one should keep in mind that oracle calls of DMP make use of very efficient MATLAB implementation of matrix-vector multiplication, while SMP relies upon much less efficient (with respect to, e.g., C language) implementation of long DO loops.

3. The advantages, if any, of SMP as compared to DMP are more significant in the case of uniform fit than in the case of  $\ell_2$  fit. This phenomenon is quite natural: as we have already explained, in the case of  $\ell_2$  fit the methods are applied to different GBSP reformulations of  $(P_2)$ , and the reformulation DMP works with is easier than the one processed by SMP.

*The results, II* In order to get impression of what happens when the matrix  $A$  in  $(P_p)$  is too large to be stored in RAM, we carried out two experiments where the goal was to solve the  $\ell_1$  minimization problem with uniform and with  $\ell_2$  fits and fully dense  $(m = 32,000) \times (n = 64,000)$  matrix  $A$  given by a simple analytical expression. This expression allows to compute a column/a row of  $A$  with a given index in  $O(m)$ , resp.,  $O(n)$  operations. Matrix  $A = A_p$  was normalized to have  $\|A\|_{1,p} = 1$ . While the sizes of  $A$  make it impossible to store the matrix in the RAM of the computer we used for the experiments, we still can multiply vectors by  $A$  and  $A^T$  by computing all necessary columns and rows, and thus can run DMP and SMP. In our related experiments, we generated at random a sparse (64 nonzeros) “true” signal  $x_* \in \mathbf{R}^{64,000}$  with  $\|x_*\|_1 = 1$ , computed  $y = Ax + \xi$ ,  $\xi$ ,  $\|\xi\|_p = \delta = 0.005$ , being observation noise, and ran DMP

**Table 3** Experiments with dense  $32,000 \times 64,000$  matrix  $A$

Method	$p$	Steps	Calls	FCalls	CPU (s)	$\ A\hat{x} - b\ _p$	$\ \hat{x} - x_*\ _r$		
							$r = 1$ (%)	$r = 2$ (%)	$r = \infty$ (%)
DMP(C00E)	$\infty$	30	71	71	7,564	0.16018	1.406 (141)	0.143 (89)	0.041 (79)
DMP(B25P)	$\infty$	31	67	67	7,363	0.15975	1.361 (136)	0.136 (85)	0.035 (69)
SMP (C00E)	$\infty$	7,501	22,141	25.9	5,352	0.00744	0.048 (5)	0.005 (3)	0.002 (4)
DMP (C00E)	2	29	67	67	7,471	0.03653	1.455 (146)	0.135 (84)	0.035 (68)
DMP (L75P)	2	30	67	67	7,536	0.02480	0.976 (98)	0.093 (58)	0.022 (42)
SMP (C00E)	2	2,602	7,749	8.5	2,350	0.00715	0.264 (26)	0.021 (13)	0.004 (7)

Percents:  $\|\hat{x} - x_*\|/\|x_*\|$

and SMP in order to find an  $\epsilon$ -solution  $x_\epsilon$ ,  $\epsilon = 0.0025$ , to the resulting problem  $(P_p)$ ; in particular, we should have  $\|x_\epsilon\|_1 \leq \|x_*\|_1 = 1$  and  $\|Ax_\epsilon - b\| \leq \delta + \epsilon = 0.0075$ . In every experiment, each of the methods was allowed to run at most 7,200 s.<sup>8</sup> The results are as follows.

- In the allowed 7,200s, the deterministic algorithms on every one of the two test problems ( $p = 2$  and  $p = \infty$ ) was able to carry out just about 30 steps with the total of about 67 computations of  $F(\cdot)$ ; this is by far not enough to get meaningful results, see Table 3. In contrast to this, the numbers of steps and randomized computations of  $F$  carried out by the randomized algorithm in the same 7,200 s was in the range of tens of thousands, which was enough to fully achieve the required accuracy for both  $p = \infty$  and  $p = 2$ .
- While the quality of approximation of  $x_*$  by the solution yielded by DMP is basically nonexistent, the SMP produced fairly reasonable approximations of  $x_*$ , see Table 3.

In our opinion, the preliminary numerical results we have reported suggest that “acceleration via randomization” possesses a significant practical potential when solving extremely large-scale convex programs of appropriate structure.

**Acknowledgments** The authors wish to express their gratitude to the Associate Editor and anonymous referees for their constructive criticism which led to substantial improvements of the paper.

## A Appendix

### A.1 Representing a vector from $\Delta_{n,d}$ as a convex combination of extreme points

We use the notations of Sect. 2.1.2. The case of  $d = n$  is trivial, thus, let  $d < n$ . Let

$$q \in \Delta_{n,d} = \left\{ q \in \mathbf{R}_+^n : 0 \leq q_i \leq 1 \forall i, \sum_{i=1}^n q_i = d \right\}.$$

<sup>8</sup> The running time is compared with the limit of 7,200s only at the end of an iteration, thus, with termination due to CPU limit, the actual running time was larger than this limit.

To represent  $q$  as a convex combination of  $n$  extreme points of  $\Delta_{n,d}$  we act as follows:

- *Initialization:* We set  $p^0 = [1; q]$ ,  $\mu^0 = 1$ . Note that  $p^0 \in \Delta = \{p = [1; p_1; \dots; p_n] \in \Delta_{n+1,d+1}\}$ .
- *Step  $t = 1, 2, \dots$ :* Given  $p^{t-1} = [1; p_1^{t-1}; \dots; p_n^{t-1}] \in \Delta$ , we find the  $d + 1$  largest among the entries  $p_i^{t-1}$ ,  $i = 1, \dots, n$ , let their indexes be  $i_1, \dots, i_{d+1}$ , where  $p_{i_1}^{t-1} \geq p_{i_2}^{t-1} \geq \dots \geq p_{i_{d+1}}^{t-1}$ .
  - (a) It may happen that  $p_{i_\ell}^{t-1} = 1$  for  $1 \leq \ell \leq d$ ; since  $p^{t-1} \in \Delta$ ,  $r^t := p^{t-1}$  is a Boolean vector with exactly  $d + 1$  entries equal to 1, and  $q^t = [p_1^{t-1}; \dots; p_n^{t-1}]$  is an extreme point of  $\Delta_{n,d}$ . We set  $v_t = 1$ ,  $p^t = 0$  and terminate.
  - (b) When not all  $p_{i_\ell}^{t-1}$ ,  $1 \leq \ell \leq d$ , are equal to 1, we set  $v_t = \min[1 - p_{i_{d+1}}^{t-1}, p_{i_d}^{t-1}]$ , define  $r^t$  as Boolean  $(n + 1)$ -dimensional vector with  $d + 1$  entries equal to 1, the indexes of the entries being  $0, i_1, \dots, i_d$ , set  $p^t = [p^{t-1} - v_t r^t] / (1 - v_t)$ ,  $q^t = [r_1^t; \dots; r_n^t]$  (note that  $q^t$  is an extreme point of  $\Delta_{n,d}$ ) and pass to step  $t + 1$ .

Observe that the algorithm is well defined. Indeed,  $0 \leq v_t \leq 1$  by construction, and  $v_t = 1$  if and only if  $p_{i_{d+1}}^{t-1} = 0$  and  $p_{i_d}^{t-1} = 1$ , that is, when we terminate at step  $t$  according to (a). Thus,  $p^t$  is well defined at every non-termination step  $t$ . Moreover, from (b) it is immediately seen that at such a step we have  $p_0^t = 1$ ,  $0 \leq p_i^t \leq 1$  for all  $i$  and  $\sum_{i=0}^n p_i^t = d + 1$ , that is,  $p^t \in \Delta$  for all  $t$  for which  $p^t$  is well defined. Besides this, it is immediately seen that  $0/1$  entries in  $p^{t-1}$  remain intact when passing from  $p^{t-1}$  to  $p^t$ , and that the total number of these entries increases at every step of the algorithm by at least 1. The latter observation implies that the algorithm terminates in at most  $n$  steps. Finally, by construction  $p^{t-1} = (1 - v_t)p^t + v_t r^t$ , whence, denoting by  $\bar{t}$  the termination step,  $p^0$  is a convex combination of  $r^1, \dots, r^{\bar{t}}$  with coefficients  $\mu_t$  readily given by  $v_1, \dots, v_{\bar{t}}$ . Discarding in  $r^1, \dots, r^{\bar{t}}$  the entries with index 0, we get extreme points  $q^1, \dots, q^{\bar{t}}$  of  $\Delta_{n,d}$  such that  $q = \sum_{i=1}^{\bar{t}} \mu_i q^i$ . Finally, the computational effort per step clearly does not exceed  $O(1)dn$ . In fact, when  $d > \ln(n)$ , at each step we can first sort the components of  $p^{t-1}$  in nonincreasing order resulting in a complexity of  $O(1)n \ln(n)$  per iteration instead of  $O(1)dn$  complexity. That is, the total computational effort is at most  $O(1) \min\{d, \ln(n)\}n^2$ .

### A.2 $\kappa$ -regular spaces

Due to space limitations, we present here a kind of “executive summary” which is fully sufficient for our purposes. For underlying definitions and proofs, see [9]; in a slightly different form, this material can be found also in [6, 16]. Consider a finite-dimensional linear space  $E$  equipped with a norm  $\| \cdot \|$ . The pair  $(E, \| \cdot \|)$  can be assigned with a well-defined *regularity parameter*  $\kappa \geq 1$  in such a way that whenever  $\xi^1, \xi^2, \dots$  are random vectors in  $E$  which form a martingale-difference, one has

$$\mathbf{E} \left\{ \|\xi^1 + \dots + \xi^N\|^2 \right\} \leq \kappa \sum_{i=1}^N \mathbf{E} \{ \|\xi^i\|^2 \}, \quad N = 1, 2, \dots$$

In addition, if  $\sigma_i$  are positive deterministic reals such that  $\mathbf{E}_{|i-1}\{\exp\{\|\xi^i\|^2/\sigma_i^2\}\} \leq \exp\{1\}$  almost surely for all  $i$ , where  $\mathbf{E}_{|i-1}$  stands for conditional,  $\xi^1, \dots, \xi^{i-1}$  being fixed, expectation, then

$$\forall(N \geq 1, \gamma \geq 0) : \text{Prob} \left\{ \|\xi^1 + \dots + \xi^N\| > [\sqrt{2\kappa} + \sqrt{2\gamma}] \sqrt{\sum_{i=1}^N \sigma_i^2} \right\} \leq \exp\{-\gamma^2/3\}.$$

In addition,

1.  $(E, \|\cdot\|)$  is  $(\dim E)$ -regular;
2. If  $(E, \|\cdot\|)$  is  $\kappa$ -regular, so is  $(E, \|\cdot\|_Q)$ , where  $Q$  is a linear automorphism of  $E$  and  $\|x\|_Q = \|Qx\|$ ;
3.  $(\mathbf{R}^N, \|\cdot\|_2)$  is regular with  $\kappa = 1$ ;
4. Let  $E$  be the space of  $m \times n$  block-diagonal matrices  $x$ ,  $m \leq n$ , of a given block-diagonal structure, and let the norm on  $E$  be defined as  $|x|_p = \|\sigma(x)\|_p$ , where  $p \geq 2$  and  $\sigma(x) \in \mathbf{R}^m$  is the vector of singular values of  $x$ . Then  $(E, |\cdot|_p)$  is regular with  $\kappa = O(1) \min[p, \ln(m + 1)]$ . In particular, when  $n \geq 3$ ,  $(\mathbf{R}^n, \|\cdot\|_\infty)$  is  $(2 \ln(n))$ -regular (treat vectors as diagonals of diagonal matrices), while the space  $\mathbf{R}^{n \times n}$  of  $n \times n$  matrices equipped with the spectral norm (maximum singular value) is  $6 \ln(n)$ -regular;
5. If  $(E_1, \|\cdot\|_1), \dots, (E_K, \|\cdot\|_K)$  are  $\kappa$ -regular, the pair  $(E = E_1 \times \dots \times E_K, \|[x^1; \dots; x^K]\| = \sqrt{\sum_{k=1}^K \|x^k\|^2})$  is  $2(\kappa + 1)$ -regular.

### A.3 Proof of Proposition 4

**1<sup>0</sup>.** Let us denote

$$\varphi_t = 8\Omega\mathcal{L}^2 + \sum_{\tau=1}^{t-1} \varsigma_\tau, \quad \varsigma_t = 3 \left[ \|F(\zeta_t) - F(w_t)\|_*^2 + \|F(\eta_t) - F(z_t)\|_*^2 \right] \quad (84)$$

(cf. (48)). Let us show that under the premise of Proposition 4

$$\forall \Lambda \geq 0 : \text{Prob} \left\{ \varphi_t \geq O(1) \left[ \Omega\mathcal{L}^2 + \frac{M_*^2 t}{k} \chi_*^2(k, \Lambda) \right] \right\} \leq \exp\{-\Lambda t\}, \quad (85)$$

where  $O(1)$  is an absolute constant factor. We use the following result (see, e.g., Theorem 2.1 (iii) of [9]): let  $\xi^i, \dots, \xi^k$  be  $k$  independent vectors from  $E$  with  $\|\xi^i\|_* \leq \sigma$  and  $\mathbf{E}\{\xi^i\} = 0$ , where the norm  $\|\cdot\|_*$  is  $\kappa$ -regular,  $\kappa \geq 1$ . Then for any  $u \geq 0$

$$\text{Prob} \left\{ \left\| \sum_{i=1}^k \xi^i \right\|_* \geq \left[ \sqrt{2\kappa} + u\sqrt{2} \right] \sigma \sqrt{k} \right\} \leq \exp\{-u^2/2\}. \quad (86)$$

When rewriting the above bound for  $\xi^i = F(\zeta^i) - F(w)$  and  $\xi^i = F(\eta^i) - F(z)$  and taking into account that  $\|\xi^i\|_* \leq M_*$  we obtain

$$\forall u \geq 0 : \text{Prob} \left\{ \left\| \sum_{i=1}^k \xi^i \right\|_*^2 \geq M_*^2 k (\sqrt{2\kappa} + \sqrt{2}u)^2 \right\} \leq \exp\{-u^2/2\}. \tag{87}$$

So, if we denote  $\text{Prob}_t$  conditional probability over  $\zeta_1, \eta_1, \dots, \zeta_{t-1}, \eta_{t-1}$  being fixed, we get

$$\forall u \geq 0 : \text{Prob}_t \left\{ \varsigma_t \geq \frac{24M_*^2}{k} (\kappa + u) \right\} \leq 2 \exp\{-u/2\}. \tag{88}$$

When setting  $\nu_t = \frac{\varsigma_t k}{24M_*^2}$ , we have for the conditional expectation  $\mathbf{E}_t$  over  $\zeta_1, \eta_1, \dots, \zeta_{t-1}, \eta_{t-1}$  being fixed and  $0 \leq \alpha < 1$

$$\begin{aligned} \mathbf{E}_t \left\{ \exp \left\{ \frac{\alpha}{2} \nu_t \right\} \right\} &\leq e^{\frac{\alpha\kappa}{2}} + \frac{\alpha}{2} \int_{\kappa}^{\infty} e^{\frac{\alpha u}{2}} \text{Prob}_t\{\nu_t \geq u\} du \\ &\leq e^{\frac{\alpha\kappa}{2}} + \alpha \int_{\kappa}^{\infty} \exp \left\{ -\frac{(1-\alpha)u}{2} \right\} du = \frac{1+\alpha}{1-\alpha} \exp \left\{ \frac{\alpha\kappa}{2} \right\} \end{aligned}$$

When choosing  $\alpha_* = \frac{\exp(1)-1}{\exp(1)+1}$  we get  $\mathbf{E}_t \left\{ \exp\left\{\frac{\alpha_* \nu_t}{2}\right\} \right\} \leq \exp\left\{\frac{\alpha_* \kappa}{2} + 1\right\}$ , so that

$$\begin{aligned} \mathbf{E} \left\{ \exp \left\{ \sum_{\tau=1}^t \frac{\alpha_* \nu_{\tau}}{2} \right\} \right\} &= \mathbf{E} \left\{ \mathbf{E}_t \left\{ \exp \left\{ \sum_{\tau=1}^{t-1} \frac{\alpha_* \nu_{\tau}}{2} \right\} \exp \left\{ \frac{\alpha_* \nu_t}{2} \right\} \right\} \right\} \\ &= \mathbf{E} \left\{ \exp \left\{ \sum_{\tau=1}^{t-1} \frac{\alpha_* \nu_{\tau}}{2} \right\} \mathbf{E}_t \left\{ \exp \left\{ \frac{\alpha_* \nu_t}{2} \right\} \right\} \right\} \\ &\leq \exp \left\{ t \left( \frac{\alpha_* \kappa}{2} + 1 \right) \right\} \end{aligned}$$

Hence, when applying the Tchebychev inequality we find

$$\forall \Lambda \geq 0 : \text{Prob} \left\{ \sum_{\tau=1}^t \nu_{\tau} \geq t \left( \kappa + \frac{2}{\alpha_*} (1 + \Lambda) \right) \right\} \leq \exp\{-\Lambda t\}. \tag{89}$$

When recalling that  $\varsigma_t \leq 6M_*^2$ , we conclude that

$$\begin{aligned} \forall \Lambda \geq 0 : \\ \text{Prob} \left\{ \sum_{\tau=1}^{t-1} \varsigma_{\tau} \geq \min \left[ 6M_*^2 t, \frac{24M_*^2 t}{k} \left( \kappa + \frac{2}{\alpha_*} (1 + \Lambda) \right) \right] \right\} &\leq \exp\{-\Lambda t\}. \tag{90} \end{aligned}$$

Since  $\kappa \geq 1$ ,  $\kappa + \frac{2}{\alpha_*} (1 + \Lambda) \leq O(1)\kappa_*^2(k, \Lambda)$ , and we arrive at (85).



2<sup>0</sup>. We have

$$\forall \lambda \geq 0 : \text{Prob} \left\{ \frac{R_t}{t} \geq O(1)F_*\sqrt{\frac{\Omega\lambda}{kt}} \right\} \leq e^{-\lambda}. \tag{91}$$

Indeed, since  $\mathcal{A}$  is skew-symmetric, i.e.  $\langle \mathcal{A}z, z \rangle = 0$ ,

$$\begin{aligned} r_t &= \langle F(\zeta_t), \zeta_t - w_t \rangle = \langle a + \mathcal{A}\zeta_t, \zeta_t - w_t \rangle = \langle a + \mathcal{A}w_t, \zeta_t - w_t \rangle \\ &= \langle F(w_t), \zeta_t - w_t \rangle. \end{aligned}$$

Let  $\zeta_t^i$  be the  $i$ th sample drawn when evaluating  $\zeta_t$ . We conclude that

$$\begin{aligned} \frac{R_t}{t} &= \frac{1}{t} \sum_{\tau=1}^t r_t = \frac{1}{t} \sum_{\tau=1}^t \langle F(w_\tau), \zeta_\tau - w_\tau \rangle = \frac{1}{t} \sum_{\tau=1}^t \left\langle F(w_\tau), \frac{1}{k} \sum_{i=1}^k \zeta_\tau^i - w_\tau \right\rangle \\ &= \frac{1}{tk} \sum_{\tau=1}^t \sum_{i=1}^k \langle F(w_\tau), \zeta_\tau^i - w_\tau \rangle = \frac{1}{tk} \sum_{\tau=1}^t \sum_{i=1}^k \xi_\tau^i, \end{aligned}$$

where  $\xi_\tau^i := \langle F(w_\tau), \zeta_\tau^i - w_\tau \rangle$  is a scalar martingale-difference with  $|\xi_\tau^i| \leq 2\mathcal{R}F_* \leq 2\Theta F_*$  (cf. (21)). Then by the Azuma-Hoeffding inequality [1],

$$\forall \lambda \geq 0 : \text{Prob} \left\{ \frac{R_t}{t} \geq 2\Theta F_*\sqrt{\frac{2\lambda}{kt}} \right\} \leq e^{-\lambda}, \tag{92}$$

which implies (91). We are done—when substituting the bounds (85) and (91) into (47) we get

$$\text{Prob} \left\{ \epsilon_{\text{sad}}(z^t) \geq O(1) \left[ \frac{\Omega\mathcal{L}}{t} + M_*\kappa_*(k, \Lambda)\sqrt{\frac{\Omega}{kt}} + \Theta F_*\sqrt{\frac{\lambda}{kt}} \right] \right\} \leq e^{-\Lambda t} + e^{-\lambda},$$

which is (56) [recall that  $\Theta = \sqrt{2\Omega}$  and  $M_* \leq 2\Theta\mathcal{L}$ , see (23)]. □

### References

1. Azuma, K.: Weighted sums of certain dependent random variables. *Tökuku Math. J.* **19**, 357–367 (1967)
2. Candès, E.J., Tao, T.: Decoding by linear programming. *IEEE Trans. Inform. Theory* **51**, 4203–4215 (2006)
3. Candès, E.J. : Compressive sampling. In: Sanz-Solé, M., Soria, J., Varona, J.L., Verdera, J. (eds.) *International Congress of Mathematicians, Madrid 2006*, vol. III, pp. 1437–1452. *European Mathematical Society Publishing House, Zurich* (2006)
4. Dalalyan, A.S., Juditsky, A., Spokoiny, V.: A new algorithm for estimating the effective dimension-reduction subspace. *J. Mach. Learn. Res.* **9**, 1647–1678 (2008)
5. Donoho, D., Huo, X.: Uncertainty principles and ideal atomic decomposition. *IEEE Trans. Inform. Theory* **47**(7), 2845–2862 (2001)

6. Dümbgen, L., van de Geer, S., Verhaar, M., Wellner, J.: Nemirovski's inequalities revisited. *Am. Math. Mon.* **117**(2), 138–160 (2010)
7. Grigoriadis, M.D., Khachiyan, I.G.: A sublinear-time randomized approximation algorithm for matrix games. *Oper. Res. Lett.* **18**, 53–58 (1995)
8. Hoeffding, W.: Probability inequalities for sums of bounded random variables. *J. Am. Stat. Assoc.* **58**(301), 13–30 (1963)
9. Juditsky, A., Nemirovski, A.: Large deviations of vector-valued martingales in 2-smooth normed spaces (2008) E-print: <http://www2.isye.gatech.edu/~nemirovs/LargeDevSubmitted.pdf>
10. Juditsky, A., Nemirovski, A., Tauvel, C.: Solving variational inequalities with stochastic mirror prox algorithm. *Stoch. Syst.* **1**(1), 17–58 (2011)
11. Lemaréchal, C., Nemirovski, A., Nesterov, Y.: New variants of bundle methods. *Math. Program.* **69**(1), 111–148 (1995)
12. Nemirovskii, A., Yudin, D.: Efficient methods for large-scale convex problems. *Ekonomika i Matematicheskie Metody* (in Russian), **15**(1), 135–152 (1979)
13. Nemirovski, A.: Prox-method with rate of convergence  $O(1/t)$  for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM J. Optim.* **15**, 229–251 (2004)
14. Nemirovski, A., Juditsky, A., Lan, G., Shapiro, A.: Stochastic approximation approach to stochastic programming. *SIAM J. Optim.* **19**(4), 1574–1609 (2009)
15. Nesterov, Y.: Smooth minimization of non-smooth functions—CORE Discussion Paper 2003/12, February 2003. *Math. Progr.* **103**, 127–152 (2005)
16. Pinelis, I.: Optimum bounds for the distributions of martingales in banach spaces. *Ann. Probab.* **22**(4), 1679–1706 (1994)