



HAL
open science

Solving variational inequalities with Stochastic Mirror-Prox algorithm

Anatoli B. Juditsky, Arkadii S. Nemirovskii, Claire Tauvel

► **To cite this version:**

Anatoli B. Juditsky, Arkadii S. Nemirovskii, Claire Tauvel. Solving variational inequalities with Stochastic Mirror-Prox algorithm. 2008. hal-00318043v1

HAL Id: hal-00318043

<https://hal.science/hal-00318043v1>

Preprint submitted on 3 Sep 2008 (v1), last revised 29 May 2011 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

SOLVING VARIATIONAL INEQUALITIES WITH STOCHASTIC MIRROR-PROX ALGORITHM

ANATOLI JUDITSKY*, ARKADI NEMIROVSKI†, AND CLAIRE TAUVEL‡

September 4, 2008

Abstract. In this paper we consider iterative methods for *stochastic variational inequalities* (s.v.i.) with monotone operators. Our basic assumption is that the operator possesses both smooth and nonsmooth components. Further, only noisy observations of the problem data are available. We develop a novel Stochastic Mirror-Prox (SMP) algorithm for solving s.v.i. and show that with the convenient stepsize strategy it attains the optimal rates of convergence with respect to the problem parameters. We apply the SMP algorithm to Stochastic composite minimization and describe particular applications to Stochastic Semidefinite Feasibility problem and Eigenvalue minimization.

Key words. Nash variational inequalities, stochastic convex-concave saddle-point problem, large scale stochastic approximation, reduced complexity algorithms for convex optimization

AMS subject classifications. 90C15, 65K10, 90C47

1. Introduction. Let Z be a convex compact set in Euclidean space \mathcal{E} with inner product $\langle \cdot, \cdot \rangle$, $\| \cdot \|$ be a norm on E (not necessarily the one associated with the inner product), and $F : Z \rightarrow \mathcal{E}$ be a monotone mapping:

$$(1.1) \quad \forall(z, z' \in Z) : \langle F(z) - F(z'), z - z' \rangle \geq 0$$

We are interested to approximate a solution to the variational inequality (v.i.)

$$(1.2) \quad \text{find } z_* \in Z : \langle F(z), z_* - z \rangle \leq 0 \quad \forall z \in Z$$

associated with Z, F . Note that since F is monotone on Z , the condition in (1.2) is implied by $\langle F(z_*), z - z_* \rangle \geq 0$ for all $z \in Z$, which is the standard definition of a (strong) solution to the v.i. associated with Z, F . The inverse – a solution to v.i. as defined by (1.2) (a “weak” solution) is a strong solution as well – also is true, provided, e.g., that F is continuous. An advantage of the concept of weak solution is that such a solution always exists under our assumptions (F is well defined and monotone on a convex compact set Z).

We quantify the inaccuracy of a candidate solution $z \in Z$ by the error

$$(1.3) \quad \text{Err}_{\text{vi}}(z) := \max_{u \in Z} \langle F(u), z - u \rangle;$$

note that this error is always ≥ 0 and equals zero iff z is a solution to (1.2).

In what follows we impose on F , aside of the monotonicity, the requirement

$$(1.4) \quad \forall(z, z' \in Z) : \|F(z) - F(z')\|_* \leq L\|z - z'\| + M$$

with some known constants $L \geq 0, M \geq 0$. From now on,

$$(1.5) \quad \|\xi\|_* = \max_{z: \|z\| \leq 1} \langle \xi, z \rangle$$

*LJK, Université J. Fourier, B.P. 53, 38041 Grenoble Cedex 9, France, anatoli.juditsky@imag.fr

†Georgia Institute of Technology, Atlanta, Georgia 30332, USA, nemirovs@isye.gatech.edu, research of this author was partly supported by the NSF award DMI-0619977.

‡LJK, Université J. Fourier, B.P. 53, 38041 Grenoble Cedex 9, France, claire.tauvel@imag.fr

is the norm conjugate to $\|\cdot\|$.

We are interested in the case where (1.2) is solved by an iterative algorithm based on a *stochastic oracle representation* of the operator $F(\cdot)$. Specifically, when solving the problem, the algorithm acquires information on F via subsequent calls to a black box (“stochastic oracle”, SO). At i -th call, $i = 0, 1, \dots$, the oracle gets as input a search point $z_i \in Z$ (this point is generated by the algorithm on the basis of the information accumulated so far) and returns the vector $\Xi(z_i, \zeta_i)$, where $\{\zeta_i \in \mathbf{R}^N\}_{i=1}^\infty$ is a sequence of i.i.d. (and independent of the queries of the algorithm) random variables. We suppose that the Borel function $\Xi(z, \zeta)$ is such that

$$(1.6) \quad \forall z \in Z : \mathbf{E} \{\Xi(z, \zeta_1)\} = F(z), \quad \mathbf{E} \{\|\Xi(z, \zeta_i) - F(z)\|_*^2\} \leq N^2.$$

We call a monotone v.i. (1.1), augmented by a stochastic oracle (SO), a *stochastic monotone v.i.* (s.v.i.).

To motivate our goal, let us start with known results [5] on the limits of performance of iterative algorithms for solving large-scale stochastic v.i.’s. To “normalize” the situation, assume that Z is the unit Euclidean ball in $\mathcal{E} = \mathbf{R}^n$ and that n is large. In this case, the rate of convergence of a whatever algorithm for solving v.i.’s cannot be better than $O(1) \left[\frac{L}{t} + \frac{M+N}{\sqrt{t}} \right]$. In other words, for a properly chosen positive absolute constant C , for every number of steps t , all large enough values of n and any algorithm \mathcal{B} for solving s.v.i.’s on the unit ball of \mathbf{R}^n , one can point out a monotone s.v.i. satisfying (1.4), (1.6) and such that the expected error of the approximate solution \tilde{z}_t generated by \mathcal{B} after t steps, applied to such s.v.i., is at least $c \left[\frac{L}{t} + \frac{M+N}{\sqrt{t}} \right]$ for some $c > 0$. To the best of our knowledge, no one of existing algorithms allows to achieve, uniformly in the dimension, this convergence rate. In fact, the “best approximations” available are given by Robust Stochastic Approximation (see [3] and references therein) with the guaranteed rate of convergence $O(1) \frac{L+M+N}{\sqrt{t}}$ and extra-gradient-type algorithms for solving deterministic monotone v.i.’s with Lipschitz continuous operators (see [6, 9, 10, 11]), which attains the accuracy $O(1) \frac{L}{t}$ in the case of $M = N = 0$ or $O(1) \frac{M}{\sqrt{t}}$ when $L = N = 0$.

The goal of this paper is to demonstrate that a specific *Mirror-Prox* algorithm [6] for solving monotone v.i.’s with Lipschitz continuous operators can be extended onto monotone s.v.i.’s to yield, uniformly in the dimension, the optimal rate of convergence $O(1) \left[\frac{L}{t} + \frac{M+N}{\sqrt{t}} \right]$. We present the corresponding extension and investigate it in details: we show how the algorithm can be “tuned” to the geometry of the s.v.i. in question, derive bounds for the probability of large deviations of the resulting error, etc. We also present a number of applications where the specific structure of the rate of convergence indeed “makes a difference”.

The main body of the paper is organized as follows: in Section 2, we describe several special cases of monotone v.i.’s we are especially interested in (convex Nash equilibria, convex-concave saddle point problems, convex minimization). We single out these special cases since here one can define a useful “functional” counterpart $\text{Err}_N(\cdot)$ of the just defined error $\text{Err}_{\text{vi}}(\cdot)$; both Err_N and Err_{vi} will participate in our subsequent efficiency estimates. Our main development – the *Stochastic Mirror Prox* (SMP) algorithm – is presented in Section 3. Some general results about the performance of the SMP are presented in Section 3.2. Then in Section 4 we present SMP for Stochastic composite minimization and discuss its applications to Stochastic Semidefinite Feasibility problem and Eigenvalue minimization. All technical proofs are collected in the appendix.

Notations. In the sequel, lowercase Latin letters denote vectors (and sometimes matrices). Script capital letters, like \mathcal{E} , \mathcal{Y} , denote Euclidean spaces; the inner product in such a space, say, \mathcal{E} , is denoted by $\langle \cdot, \cdot \rangle_{\mathcal{E}}$ (or merely $\langle \cdot, \cdot \rangle$, when the corresponding space is clear from the context). Linear mappings from one Euclidean space to another, say, from \mathcal{E} to \mathcal{F} , are denoted by boldface capitals like \mathbf{A} (there are also some reserved boldface capitals, like \mathbf{E} for expectation, \mathbf{R}^k for the k -dimensional coordinate space, and \mathbf{S}^k for the space of $k \times k$ symmetric matrices). \mathbf{A}^* stands for the conjugate to mapping \mathbf{A} : if $\mathbf{A} : \mathcal{E} \rightarrow \mathcal{F}$, then $\mathbf{A}^* : \mathcal{F} \rightarrow \mathcal{E}$ is given by the identity $\langle f, \mathbf{A}e \rangle_{\mathcal{F}} = \langle \mathbf{A}^* f, e \rangle_{\mathcal{E}}$ for $f \in \mathcal{F}, e \in \mathcal{E}$. When both the origin and the destination space of a linear map, like \mathbf{A} , are the standard coordinate spaces, the map is identified with its matrix A , and \mathbf{A}^* is identified with A^T . For a norm $\| \cdot \|$ on \mathcal{E} , $\| \cdot \|_*$ stands for the conjugate norm, see (1.5).

For Euclidean spaces $\mathcal{E}_1, \dots, \mathcal{E}_m$, $\mathcal{E} = \mathcal{E}_1 \times \dots \times \mathcal{E}_m$ denotes their Euclidean direct product, so that a vector from \mathcal{E} is a collection $u = [u_1; \dots; u_m]$ (“MATLAB notation”) of vectors $u_\ell \in \mathcal{E}_\ell$, and $\langle u, v \rangle_{\mathcal{E}} = \sum_{\ell} \langle u_\ell, v_\ell \rangle_{\mathcal{E}_\ell}$. Sometimes we allow ourselves to write (u_1, \dots, u_m) instead of $[u_1; \dots; u_m]$.

2. Preliminaries.

2.1. Nash v.i.’s and functional error. In the sequel, we shall be especially interested in a special case of v.i. (1.2) – in a *Nash v.i.* coming from a *convex Nash Equilibrium* problem, and in the associated *functional* error measure. The Nash Equilibrium problem can be described as follows: there are m players, i -th of them choosing a point z_i from a given set Z_i . The loss of i -th player is a given function $\phi_i(z)$ of the collection $z = (z_1, \dots, z_m) \in Z = Z_1 \times \dots \times Z_m$ of player’s choices. With slight abuse of notation, we use for $\phi_i(z)$ also the notation $\phi_i(z_i, z^i)$, where z^i is the collection of choices of all but the i -th players. Players are interested to minimize their losses, and Nash equilibrium \hat{z} is a point from Z such that for every i the function $\phi_i(z_i, \hat{z}^i)$ attains its minimum in $z_i \in Z_i$ at $z_i = \hat{z}_i$ (so that in the state \hat{z} no player has an incentive to change his choice, provided that the other players stick to their choices).

We call a Nash equilibrium problem *convex*, if for every i Z_i is a compact convex set, $\phi_i(z_i, z^i)$ is a Lipschitz continuous function convex in z_i and concave in z^i , and the function $\Phi(z) = \sum_{i=1}^m \phi_i(z)$ is convex. It is well known (see, e.g., [8]) that setting

$$F(z) = [F^1(z); \dots; F^m(z)], \quad F^i(z) \in \partial_{z_i} \phi_i(z_i, z^i), \quad i = 1, \dots, m$$

where $\partial_{z_i} \phi_i(z_i, z^i)$ is the subdifferential of the convex function $\phi_i(\cdot, z^i)$ at a point z_i , we get a monotone operator such that the solutions to the corresponding v.i. (1.2) are exactly the Nash equilibria. Note that since ϕ_i are Lipschitz continuous, the associated operator F can be chosen to be bounded. For this v.i. one can consider, along with the v.i.-accuracy measure $\text{Err}_{\text{v.i.}}(z)$, the *functional* error measure

$$\text{Err}_{\text{N}}(z) = \sum_{i=1}^m \left[\phi_i(z) - \min_{w_i \in Z_i} \phi_i(w_i, z^i) \right]$$

This accuracy measure admits a transparent justification: this is the sum, over the players, of the incentives for a player to change his choice given that other players stick to their choices.

Special cases: saddle points and minimization. An important by its own right particular case of Nash Equilibrium problem is an *antagonistic 2-person game*, where

$m = 2$ and $\Phi(z) \equiv 0$ (i.e., $\phi_2(z) \equiv -\phi_1(z)$). The convex case of this problem corresponds to the situation when $\phi(z_1, z_2) \equiv \phi_1(z_1, z_2)$ is a Lipschitz continuous function which is convex in $z_1 \in Z_1$ and concave in $z_2 \in Z_2$, the Nash equilibria are exactly the saddle points (min in z_1 , max in z_2) of ϕ on $Z_1 \times Z_2$, and the functional error becomes

$$\text{Err}_N(z_1, z_2) = \max_{(u_1, u_2) \in Z} [\phi(z_1, u_1) - \phi(u_1, z_2)].$$

Recall that the convex-concave saddle point problem $\min_{z_1 \in Z_1} \max_{z_2 \in Z_2} \phi(z_1, z_2)$ gives rise to the “primal-dual” pair of convex optimization problems

$$(P) : \min_{z_1 \in Z_1} \bar{\phi}(z_1), \quad (D) : \max_{z_2 \in Z_2} \underline{\phi}(z_2),$$

where

$$\bar{\phi}(z_1) = \max_{z_2 \in Z_2} \phi(z_1, z_2), \quad \underline{\phi}(z_2) = \min_{z_1 \in Z_1} \phi(z_1, z_2).$$

The optimal values $\text{Opt}(P)$ and $\text{Opt}(D)$ in these problems are equal, the set of saddle points of ϕ (i.e., the set of Nash equilibria of the underlying convex Nash problem) is exactly the direct product of the optimal sets of (P) and (D) , and $\text{Err}_N(z_1, z_2)$ is nothing but the sum of non-optimality of z_1, z_2 considered as approximate solutions to respective optimization problems:

$$\text{Err}_N(z_1, z_2) = [\bar{\phi}(z_1) - \text{Opt}(P)] + [\text{Opt}(D) - \underline{\phi}(z_2)].$$

Finally, the “trivial” case $m = 1$ of the convex Nash Equilibrium is the problem of minimizing a Lipschitz continuous convex function $\phi(z) = \phi_1(z_1)$ over the convex compact set $Z = Z_1$. In this case, the functional error becomes the usual residual in terms of the objective:

$$\text{Err}_N(z) = \phi(z) - \min_Z \phi.$$

In the sequel, we refer to the v.i. (1.2) coming from a convex Nash Equilibrium problem as *Nash v.i.*, and to the two just outlined particular cases of the Nash v.i. as the *Saddle Point* and the *Minimization v.i.*, respectively. It is easy to verify that in the Saddle Point/Minimization case the functional error $\text{Err}_N(z)$ is $\geq \text{Err}_{\text{vi}}(z)$; this is not necessary so for a general Nash v.i.

2.2. Prox-mapping. We once for ever fix a norm $\|\cdot\|$ on \mathcal{E} ; $\|\cdot\|_*$ stands for the conjugate norm, see (1.5). A *distance-generating function* for Z is, by definition, a continuous convex function $\omega(\cdot) : Z \rightarrow \mathbf{R}$ such that

1. if Z^o be the set of all points $z \in Z$ such that the subdifferential $\partial\omega(z)$ of $\omega(\cdot)$ at z is nonempty, then the subdifferential of ω admits a continuous selection on Z^o : there exists a continuous on Z^o vector-valued function $\omega'(z)$ such that $\omega'(z) \in \partial\omega(z)$ for all $z \in Z^o$;
2. for certain $\alpha > 0$, $\omega(\cdot)$ is strongly convex, modulus α , w.r.t. the norm $\|\cdot\|$:

$$(2.1) \quad \forall(z, z' \in Z^o) : \langle \omega'(z) - \omega'(z'), z - z' \rangle \geq \alpha \|z - z'\|^2.$$

In the sequel, we fix a distance-generating function $\omega(\cdot)$ for Z and assume that $\omega(\cdot)$ and Z “fit” each other, meaning that one can easily solve problems of the form

$$(2.2) \quad \min_{z \in Z} [\omega(z) + \langle e, z \rangle], \quad e \in \mathcal{E}.$$

The *prox-function* associated with the distance-generating function ω is defined as

$$V(z, u) = \omega(u) - \omega(z) - \langle \omega'(z), u - z \rangle : Z^o \times Z \rightarrow \mathbf{R}^+.$$

We set

$$(2.3) \quad \begin{array}{ll} (a) & \Theta(z) = \max_{u \in Z} V(z, u) \quad [z \in Z^o]; \\ (c) & \Theta = \Theta(z_c); \end{array} \quad \begin{array}{ll} (b) & z_c = \operatorname{argmin}_Z \omega(z); \\ (d) & \Omega = \sqrt{2\Theta/\alpha}. \end{array}$$

Note that z_c is well defined (since Z is a convex compact set and $\omega(\cdot)$ is continuous and strongly convex on Z) and belongs to Z^o (since $0 \in \partial\omega(z_c)$). Note also that due to the strong convexity of ω and the origin of z_c we have

$$(2.4) \quad \forall (u \in Z) : \frac{\alpha}{2} \|u - z_c\|^2 \leq \Theta \leq \max_{z \in Z} \omega(z) - \omega(z_c);$$

in particular we see that

$$(2.5) \quad Z \subset \{z : \|z - z_c\| \leq \Omega\}.$$

Prox-mapping. Given $z \in Z^o$, we associate with this point and $\omega(\cdot)$ the *prox-mapping*

$$P(z, \xi) = \operatorname{argmin}_{u \in Z} \{\omega(u) + \langle \xi - \omega'(z), u \rangle\} \equiv \operatorname{argmin}_{u \in Z} \{V(z, u) + \langle \xi, u \rangle\} : \mathcal{E} \rightarrow Z^o.$$

We illustrate the just-defined notions with three basic examples.

Example 1: Euclidean setup. Here \mathcal{E} is \mathbf{R}^N with the standard inner product, $\|\cdot\|_2$ is the standard Euclidean norm on \mathbf{R}^N (so that $\|\cdot\|_* = \|\cdot\|$) and $\omega(z) = \frac{1}{2}z^T z$ (i.e., $Z^o = Z$, $\alpha = 1$). Assuming for the sake of simplicity that $0 \in Z$, $z_c = 0$, $\Omega = \max_{z \in Z} \|z\|_2$ and $\Theta = \frac{1}{2}\Omega^2$. The prox-function and the prox-mapping are given by $V(z, u) = \frac{1}{2}\|z - u\|_2^2$, $P(z, \xi) = \operatorname{argmin}_{u \in Z} \|(z - \xi) - u\|_2$.

Example 2: Simplex setup. Here \mathcal{E} is \mathbf{R}^N , $N > 1$, with the standard inner product, $\|z\| = \|z\|_1 := \sum_{j=1}^N |z_j|$ (so that $\|\xi\|_* = \max_j |\xi_j|$), Z is a closed convex subset of the standard simplex

$$\mathcal{D}_N = \{z \in \mathbf{R}^N : z \geq 0, \sum_{j=1}^N z_j = 1\}$$

containing its barycenter, and $\omega(z) = \sum_{j=1}^N z_j \ln z_j$ is the entropy. Then

$$Z^o = \{z \in Z : z > 0\} \quad \text{and} \quad \omega'(z) = [1 + \ln z_1; \dots; 1 + \ln z_N], \quad z \in Z^o.$$

It is easily seen (see, e.g., [3]) that here

$$\alpha = 1, \quad z_c = [1/N; \dots; 1/N], \quad \Theta \leq \ln(N)$$

(the latter inequality becomes equality when Z contains a vertex of \mathcal{D}_N), and thus $\Omega \leq \sqrt{2 \ln N}$. The prox-function is

$$V(z, u) = \sum_{j=1}^N u_j \ln(u_j/z_j),$$

and the prox-mapping is easy to compute when $Z = \mathcal{D}_N$:

$$(P(z, \xi))_j = \left(\sum_{i=1}^N z_i \exp\{-\xi_i\} \right)^{-1} z_j \exp\{-\xi_j\}.$$

Example 3: Spectahedron setup. This is the “matrix analogy” of the Simplex setup. Specifically, now \mathcal{E} is the space of $N \times N$ block-diagonal symmetric matrices, $N > 1$, of a given block-diagonal structure equipped with the Frobenius inner product $\langle a, b \rangle_F = \text{Tr}(ab)$ and the trace norm $|a|_1 = \sum_{i=1}^N |\lambda_i(a)|$, where $\lambda_1(a) \geq \dots \geq \lambda_N(a)$ are the eigenvalues of a symmetric $N \times N$ matrix a ; the conjugate norm $|a|_\infty$ is the usual spectral norm (the largest singular value) of a . Z is assumed to be a closed convex subset of the *spectahedron* $\mathcal{S} = \{z \in \mathcal{E} : z \succeq 0, \text{Tr}(z) = 1\}$ containing the matrix $N^{-1}I_N$. The distance-generating function is the matrix entropy

$$\omega(z) = \sum_{j=1}^N \lambda_j(z) \ln \lambda_j(z),$$

so that $Z^\circ = \{z \in Z : z \succ 0\}$ and $\Omega'(z) = \ln(z)$. This setup, similarly to the Simplex one, results in $\alpha = 1$, $z_c = N^{-1}I_N$, $\Theta = \ln N$ and $\Omega = \sqrt{2 \ln N}$ [2]. When $Z = \mathcal{S}$, it is relatively easy to compute the prox-mapping (see [2, 6]); this task reduces to the singular value decomposition of a matrix from \mathcal{E} . It should be added that the matrices from \mathcal{S} are exactly the matrices of the form

$$a = \mathcal{H}(b) \equiv (\text{Tr}(\exp\{b\}))^{-1} \exp\{b\}$$

with $b \in \mathcal{E}$. Note also that when $Z = \mathcal{S}$, the prox-mapping becomes “linear in matrix logarithm”: if $z = \mathcal{H}(a)$, then $P(z, \xi) = \mathcal{H}(a - \xi)$.

3. Stochastic Mirror-Prox algorithm.

3.1. Mirror-Prox algorithm with erroneous information. We are about to present the Mirror-Prox algorithm proposed in [6]. In contrast to the original version of the method, below we allow for errors when computing the values of F – we assume that given a point $z \in Z$, we can compute an approximation $\widehat{F}(z) \in \mathcal{E}$ of $F(z)$. The t -step Mirror-Prox algorithm as applied to (1.2) is as follows:

ALGORITHM 3.1.

1. Initialization: Choose $r_0 \in Z^\circ$ and stepsizes $\gamma_\tau > 0$, $1 \leq \tau \leq t$.
2. Step τ , $\tau = 1, 2, \dots, t$: Given $r_{\tau-1} \in Z^\circ$, set

$$(3.1) \quad \begin{cases} w_\tau &= P(r_{\tau-1}, \gamma_\tau \widehat{F}(r_{\tau-1})), \\ r_\tau &= P(r_{\tau-1}, \gamma_\tau \widehat{F}(w_\tau)) \end{cases}$$

. When $\tau < t$, loop to step $t + 1$.

3. At step t , output

$$(3.2) \quad \widehat{z}_t = \left[\sum_{\tau=1}^t \gamma_\tau \right]^{-1} \sum_{\tau=1}^t \gamma_\tau w_\tau.$$

The preliminary technical result on the outlined algorithm is as follows.

THEOREM 3.2. Consider t -step algorithm 3.1 as applied to a v.i. (1.2) with a monotone operator F satisfying (1.4). For $\tau = 1, 2, \dots$, let us set

$$\Delta_\tau = F(w_\tau) - \widehat{F}(w_\tau);$$

for z belonging to the trajectory $\{r_0, w_1, r_1, \dots, w_t, r_t\}$ of the algorithm, let

$$\epsilon_z = \|\widehat{F}(z) - F(z)\|_*,$$

and let $\{y_\tau \in Z^o\}_{\tau=0}^t$ be the sequence given by the recurrence

$$(3.3) \quad y_\tau = P(y_{\tau-1}, \gamma_\tau \Delta_\tau), \quad y_0 = r_0.$$

Assume that

$$(3.4) \quad \gamma_\tau \leq \frac{\alpha}{\sqrt{3L}},$$

Then

$$(3.5) \quad \text{Err}_{\text{vi}}(\hat{z}_t) \leq \left(\sum_{\tau=1}^t \gamma_\tau \right)^{-1} \Gamma(t),$$

where $\text{Err}_{\text{vi}}(\hat{z}_t)$ is defined in (1.3),

$$(3.6) \quad \Gamma(t) = 2\Theta(r_0) + \sum_{\tau=1}^t \frac{3\gamma_\tau^2}{2\alpha} \left[M^2 + (\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2 + \frac{\epsilon_{w_\tau}^2}{3} \right] \\ + \sum_{\tau=1}^t \langle \gamma_\tau \Delta_\tau, w_\tau - y_{\tau-1} \rangle$$

and $\Theta(\cdot)$ is defined by (2.3).

Finally, when (1.2) is a Nash v.i., one can replace $\text{Err}_{\text{vi}}(\hat{z}_t)$ in (3.5) with $\text{Err}_{\text{N}}(\hat{z}_t)$.

3.2. Main result. From now on, we focus on the case when Algorithm 3.1 solves monotone v.i. (1.2), and the corresponding monotone operator F is represented by a stochastic oracle. Specifically, at the i -th call to the SO, the input being $z \in Z$, the oracle returns the vector $\hat{F} = \Xi(z, \zeta_i)$, where $\{\zeta_i \in \mathbf{R}^N\}_{i=1}^\infty$ is a sequence of i.i.d. random variables, and $\Xi(z, \zeta) : Z \times \mathbf{R}^N \rightarrow \mathcal{E}$ is a Borel function. We refer to this specific implementation of Algorithm 3.1 as to *Stochastic Mirror Prox* (SMP) algorithm.

In the sequel, we impose on the SO in question the following assumption, slightly milder than (1.6):

Assumption I: With some $\mu \in [0, \infty)$, for all $z \in Z$ we have

$$(3.7) \quad \begin{aligned} (a) \quad & \|\mathbf{E} \{\Xi(z, \zeta_i) - F(z)\}\|_* \leq \mu \\ (b) \quad & \mathbf{E} \{\|\Xi(z, \zeta_i) - F(z)\|_*^2\} \leq M^2. \end{aligned}$$

In some cases, we augment Assumption I by the following

Assumption II: For all $z \in Z$ and all i we have

$$(3.8) \quad \mathbf{E} \left\{ \exp\{\|\Xi(z, \zeta_i) - F(z)\|_*^2/M^2\} \right\} \leq \exp\{1\}.$$

Note that Assumption II implies (3.7.b), since

$$\exp\{\mathbf{E} \{\|\Xi(z, \zeta_i) - F(z)\|_*^2/M^2\}\} \leq \mathbf{E} \left\{ \exp\{\|\Xi(z, \zeta_i) - F(z)\|_*^2/M^2\} \right\}$$

by the Jensen inequality.

REMARK 3.3. Observe that that the accuracy of Algorithm 3.1 (cf. (3.6)) depends in the same way on the “size” of perturbation $\epsilon_z = \|\hat{F}(z) - F(z)\|_*$ and the bound M of (1.4) on the variation of the non-Lipschitz component of F . This is why, to simplify

the presentation, we decided to use the same bound M for the scale of perturbation $\Xi(z, \zeta_i) - F(z)$ in (3.7), (3.8).

REMARK 3.4. From now on, we assume that the starting point r_0 in Algorithm 3.1 is the minimizer z_c of $\omega(\cdot)$ on Z . Further, to avoid unnecessarily complicated formulas (and with no harm to the efficiency estimates) we stick to the constant stepsize policy $\gamma_\tau \equiv \gamma$, $1 \leq \tau \leq t$, where t is a fixed in advance number of iterations of the algorithm. Our main result is as follows:

THEOREM 3.5. Let v.i. (1.2) with monotone operator F satisfying (1.4) be solved by t -step Algorithm 3.1 using a SO, and let the stepsizes $\gamma_\tau \equiv \gamma$, $1 \leq \tau \leq t$, satisfy $0 < \gamma \leq \frac{\alpha}{\sqrt{3}L}$, see (1.4). Then

(i) Under Assumption I, one has

$$(3.9) \quad \mathbf{E} \{ \text{Err}_{\text{vi}}(\widehat{z}_t) \} \leq K_0(t) \equiv \left[\frac{\alpha\Omega^2}{t\gamma} + \frac{21M^2\gamma}{2\alpha} \right] + 2\mu\Omega,$$

where M is the constant from (1.4) and Ω is given by (2.3).

(ii) Under Assumptions I, II, one has, in addition to (3.9), for any $\Lambda > 0$,

$$(3.10) \quad \text{Prob} \{ \text{Err}_{\text{vi}}(\widehat{z}_t) > K_0(t) + \Lambda K_1(t) \} \leq \exp\{-\Lambda^2/3\} + \exp\{-\Lambda t\},$$

where

$$K_1(t) = \frac{7M^2\gamma}{2\alpha} + \frac{2M\Omega}{\sqrt{t}}.$$

In the case of a Nash v.i., $\text{Err}_{\text{vi}}(\cdot)$ in (3.9), (3.10) can be replaced with $\text{Err}_{\text{N}}(\cdot)$.

When optimizing the bound (3.9) in γ , we get the following

COROLLARY 3.6. In the situation of Theorem 3.5, let the stepsizes $\gamma_\tau \equiv \gamma$ be chosen according to

$$(3.11) \quad \gamma = \min \left[\frac{\alpha}{\sqrt{3}L}, \frac{\alpha\Omega}{M} \sqrt{\frac{2}{21t}} \right].$$

Then under Assumption I one has

$$(3.12) \quad \mathbf{E} \{ \text{Err}_{\text{vi}}(\widehat{z}_t) \} \leq K_0^*(t) \equiv \max \left[\frac{7}{4} \frac{\Omega^2 L}{t} + 7 \frac{\Omega M}{\sqrt{t}} \right] + 2\mu\Omega, .$$

(see (2.3)). Under Assumptions I, II, one has, in addition to (3.12), for any $\Lambda > 0$,

$$(3.13) \quad \text{Prob} \{ \text{Err}_{\text{vi}}(\widehat{z}_t) > K_0^*(t) + \Lambda K_1^*(t) \} \leq \exp\{-\Lambda^2/3\} + \exp\{-\Lambda t\}$$

with

$$K_1^*(t) = \frac{7}{2} \frac{\Omega M}{\sqrt{t}}.$$

In the case of a Nash v.i., $\text{Err}_{\text{vi}}(\cdot)$ in (3.12), (3.13) can be replaced with $\text{Err}_{\text{N}}(\cdot)$.

3.3. Comparison with Robust Mirror SA Algorithm. Consider the case of a Nash s.v.i. with operator F satisfying (1.4) with $L = 0$, and let the SO be unbiased (i.e., $\mu = 0$). In this case, the bound (3.12) reads

$$(3.14) \quad \mathbf{E} \{ \text{Err}_{\text{N}}(\widehat{z}_t) \} \leq \frac{7\Omega M}{\sqrt{t}},$$

where

$$M^2 = \max \left[\sup_{z, z' \in Z} \|F(z) - F(z')\|_*^2, \sup_{z \in Z} \mathbf{E} \{ \|\Xi(z, \zeta_i) - F(z)\|_*^2 \} \right]$$

The bound (3.14) looks very much like the efficiency estimate

$$(3.15) \quad \mathbf{E} \{ \text{Err}_N(\tilde{z}_t) \} \leq O(1) \frac{\Omega \bar{M}}{\sqrt{t}}$$

(from now on, all $O(1)$'s are appropriate absolute positive constants) for the approximate solution \tilde{z}_t of the t -step *Robust Mirror SA* (RMSA) algorithm [3]¹⁾. In the latter estimate, Ω is exactly the same as in (3.14), and \bar{M} is given by

$$\bar{M}^2 = \max \left[\sup_z \|F(z)\|_*^2; \sup_{z \in Z} \mathbf{E} \{ \|\Xi(z, \zeta_i) - F(z)\|_*^2 \} \right].$$

Note that we always have $M \leq 2\bar{M}$, and typically M and \bar{M} are of the same order of magnitude; it may happen, however (think of the case when F is “almost constant”), that $M \ll \bar{M}$. Thus, the bound (3.14) never is worse, and sometimes can be much better than the SA bound (3.15). It should be added that as far as implementation is concerned, the SMP algorithm is not more complicated than the RMSA (cf. the description of Algorithm 3.1 with the description

$$\begin{aligned} r_t &= P(r_{t-1}, \hat{F}(r_{t-1})), \\ \hat{z}_t &= \left[\sum_{\tau=1}^t \gamma_\tau \right]^{-1} \sum_{\tau=1}^t \gamma_\tau r_\tau, \end{aligned}$$

of the RMSA).

The just outlined advantage of SMP as compared to the usual Stochastic Approximation is not that important, since “typically” M and \bar{M} are of the same order. We believe that the most interesting feature of the SMP algorithm is its ability to take advantage of a specific structure of a stochastic optimization problem, namely, insensitivity to the presence in the objective of large, but smooth and well-observable components.

We are about to consider several less straightforward applications of the outlined insensitivity of the SMP algorithm to smooth well-observed components in the objective.

4. Application to Stochastic Approximation: Stochastic composite minimization.

4.1. Problem description. Consider the optimization problem as follows (cf. [5]):

$$(4.1) \quad \min_{x \in X} \phi(x) := \Phi(\phi_1(x), \dots, \phi_m(x)),$$

where

¹⁾ In this reference, only the Minimization and the Saddle Point problems are considered. However, the results of [3] can be easily extended to s.v.i.'s.

1. $X \subset \mathcal{X}$ is a convex compact; the embedding space \mathcal{X} is equipped with a norm $\|\cdot\|_x$, and X - with a distance-generating function $\omega_x(x)$ with certain parameters $\alpha_x, \Theta_x, \Omega_x$ w.r.t. the norm $\|\cdot\|_x$;
2. $\phi_\ell(x) : X \rightarrow \mathcal{E}_\ell$, $1 \leq \ell \leq m$, are Lipschitz continuous mappings taking values in Euclidean spaces \mathcal{E}_ℓ equipped with norms (not necessarily the Euclidean ones) $\|\cdot\|_{(\ell)}$ with conjugates $\|\cdot\|_{(\ell,*)}$ and with closed convex cones K_ℓ . We suppose that ϕ_ℓ are K_ℓ -convex, i.e. for any $x, x' \in X$, $\lambda \in [0, 1]$,

$$\phi_\ell(\lambda x + (1 - \lambda)x') \leq_{K_\ell} \lambda \phi_\ell(x) + (1 - \lambda)\phi_\ell(x'),$$

where the notation $a \leq_K b \Leftrightarrow b \geq_K a$ means that $b - a \in K$.

In addition to these structural restrictions, we assume that for all $v, v' \in X$, $h \in \mathcal{X}$,

$$(4.2) \quad \begin{aligned} (a) \quad & \|[\phi'_\ell(v) - \phi'_\ell(v')]h\|_{(\ell)} \leq [L_x \|v - v'\|_x + M_x] \|h\|_x \\ (b) \quad & \|[\phi'_\ell(v)]h\|_{(\ell)} \leq [L_x \Omega_x + M_x] \|h\|_x \end{aligned}$$

for certain selections $\phi'_\ell(x) \in \partial^{K_\ell} \phi_\ell(x)$, $x \in X^2$) and certain nonnegative constants L_x and M_x .

3. Functions $\phi_\ell(\cdot)$ are represented by an unbiased SO. At i -th call to the oracle, $x \in X$ being the input, the oracle returns vectors $f_\ell(x, \zeta_i) \in \mathcal{E}_\ell$ and linear mappings $\mathbf{G}_\ell(x, \zeta_i)$ from \mathcal{X} to \mathcal{E}_ℓ , $1 \leq \ell \leq m$ ($\{\zeta_i\}$ are i.i.d. random vectors) such that for any $x \in X$ and $i = 1, 2, \dots$,

$$(4.3) \quad \begin{aligned} (a) \quad & \mathbf{E} \{f_\ell(x, \zeta_i)\} = \phi_\ell(x), \quad 1 \leq \ell \leq m \\ (b) \quad & \mathbf{E} \left\{ \max_{1 \leq \ell \leq m} \|f_\ell(x, \zeta_i) - \phi_\ell(x)\|_{(\ell)}^2 \right\} \leq M_x^2 \Omega_x^2; \\ (c) \quad & \mathbf{E} \{\mathbf{G}_\ell(x, \zeta_i)\} = \phi'_\ell(x), \quad 1 \leq \ell \leq m, \\ (d) \quad & \mathbf{E} \left\{ \max_{\substack{h \in \mathcal{X} \\ \|h\|_x \leq 1}} \|[\mathbf{G}_\ell(x, \zeta_i) - \phi'_\ell(x)]h\|_{(\ell)}^2 \right\} \leq M_x^2, \quad 1 \leq \ell \leq m. \end{aligned}$$

4. $\Phi(\cdot)$ is a convex function on $\mathcal{E} = \mathcal{E}_1 \times \dots \times \mathcal{E}_m$ given by the representation

$$(4.4) \quad \Phi(u_1, \dots, u_m) = \max_{y \in Y} \left\{ \sum_{\ell=1}^m \langle u_\ell, \mathbf{A}_\ell y + b_\ell \rangle_{\mathcal{E}_\ell} - \Phi_*(y) \right\}$$

for $u_\ell \in \mathcal{E}_\ell$, $1 \leq \ell \leq m$. Here

- (a) $Y \subset \mathcal{Y}$ is a convex compact set containing the origin; the embedding Euclidean space \mathcal{Y} is equipped with a norm $\|\cdot\|_y$, and Y - with a distance-generating function $\omega_y(y)$ with parameters $\alpha_y, \Theta_y, \Omega_y$ w.r.t. the norm $\|\cdot\|_y$;
- (b) The affine mappings $y \mapsto \mathbf{A}_\ell y + b_\ell : \mathcal{Y} \rightarrow \mathcal{E}_\ell$ are such that $\mathbf{A}_\ell y + b_\ell \in K_\ell^*$ for all $y \in Y$ and all ℓ ; here K_ℓ^* is the cone dual to K_ℓ ;

²⁾ For a K -convex function $\phi : X \rightarrow \mathcal{E}$ ($X \subset \mathcal{X}$ is convex, $K \subset \mathcal{E}$ is a closed convex cone) and $x \in X$, the K -subdifferential $\partial^K \phi(x)$ is comprised of all linear mappings $h \mapsto \mathbf{P}h : \mathcal{X} \rightarrow \mathcal{E}$ such that $\phi(u) \geq_K \phi(x) + \mathbf{P}(u - x)$ for all $u \in X$. When ϕ is Lipschitz continuous on X , $\partial^K \phi(x) \neq \emptyset$ for all $x \in X$; if ϕ is differentiable at $x \in \text{int } X$ (as it is the case almost everywhere on $\text{int } X$), one has $\frac{\partial \phi(x)}{\partial x} \in \partial^K \phi(x)$.

(c) $\Phi_*(y)$ is a given convex function on Y such that

$$(4.5) \quad \|\Phi'_*(y) - \Phi'_*(y')\|_{y,*} \leq L_y \|y - y'\|_y + M_y$$

for certain selection $\Phi'_*(z) \in \partial\Phi_*(y)$, $y \in Y$.

Example: Stochastic Matrix Minimax problem (SMMP). For $1 \leq \ell \leq m$, let $\mathcal{E}_\ell = \mathbf{S}^{p_\ell}$ be the space of symmetric $p_\ell \times p_\ell$ matrices equipped with the Frobenius inner product $\langle A, B \rangle_F = \text{Tr}(AB)$ and the spectral norms $|\cdot|_\infty$, and let K_ℓ be the cone $\mathbf{S}_+^{p_\ell}$ of symmetric positive semidefinite $p_\ell \times p_\ell$ matrices. Consider the problem

$$(P) \quad \min_{x \in X} \max_{1 \leq j \leq k} \lambda_{\max} \left(\sum_{\ell=1}^m P_{j\ell}^T \phi_\ell(x) P_{j\ell} \right),$$

where $P_{j\ell}$ are given $p_\ell \times q_j$ matrices, and $\lambda_{\max}(A)$ is the maximal eigenvalue of a symmetric matrix A . Observing that for a symmetric $p \times q$ matrix A one has

$$\lambda_{\max}(A) = \max_{S \in \mathcal{S}_q} \text{Tr}(AS)$$

where $\mathcal{S}_q = \{S \in \mathbf{S}_+^q : \text{Tr}(S) = 1\}$. When denoting by Y the set of all symmetric positive semidefinite block-diagonal matrices $y = \text{Diag}\{y_1, \dots, y_k\}$ with unit trace and diagonal blocks y_j of sizes $q_j \times q_j$, we can represent (P) in the form of (4.1), (4.4) with

$$\begin{aligned} \Phi(u) &:= \max_{1 \leq j \leq k} \lambda_{\max} \left(\sum_{\ell=1}^m P_{j\ell} u_\ell P_{j\ell}^T \right) \\ &= \max_{y = \text{Diag}\{y_1, \dots, y_k\} \in Y} \sum_{j=1}^k \text{Tr} \left(\sum_{\ell=1}^m P_{j\ell}^T u_\ell P_{j\ell} y_k \right) \\ &= \max_{y = \text{Diag}\{y_1, \dots, y_k\} \in Y} \sum_{\ell=1}^m \text{Tr} \left(u_\ell \left[\sum_{j=1}^k P_{j\ell}^T y_k P_{j\ell} \right] \right) \\ &= \max_{y = \text{Diag}\{y_1, \dots, y_k\} \in Y} \sum_{\ell=1}^m \langle u_\ell, \mathbf{A}_\ell y \rangle_F \end{aligned}$$

(we put $\mathbf{A}_\ell y = \sum_{j=1}^k P_{j\ell}^T y_j P_{j\ell}$). The set Y is the spectahedron in the space \mathbf{S}^q of symmetric block-diagonal matrices with k diagonal blocks of the sizes $q_j \times q_j$, $1 \leq j \leq k$. When equipping Y with the spectahedron setup, we get $\alpha_y = 1$, $\Theta_y = \ln(\sum_{j=1}^k q_j)$ and $\Omega_y = \sqrt{2 \ln(\sum_{j=1}^k q_j)}$, see Section 2.2.

Observe that in the simplest case of $k = m$, $p_j = q_j$, $1 \leq j \leq m$ and $P_{j\ell}$ equal to I_p for $j = \ell$ and to 0 otherwise, the SMMP problem becomes

$$(4.6) \quad \min_{x \in X} \left[\max_{1 \leq \ell \leq m} \lambda_{\max}(\phi_\ell(x)) \right].$$

If, in addition, $p_j = q_j = 1$ for all j , we arrive at the usual (“scalar”) minimax problem

$$(4.7) \quad \min_{x \in X} \left[\max_{1 \leq \ell \leq m} \phi_\ell(x) \right]$$

with convex real-valued functions ϕ_ℓ .

Observe that in the case of (4.4), the optimization problem (4.1) is nothing but the primal problem associated with the saddle point problem

$$(4.8) \quad \min_{x \in X} \max_{y \in Y} \left[\phi(x, y) = \sum_{\ell=1}^m \langle \phi_\ell(x), \mathbf{A}_\ell y + b_\ell \rangle_{\mathcal{E}_\ell} - \Phi_*(y) \right]$$

and the cost function in the latter problem is Lipschitz continuous and convex-concave due to the K_ℓ -convexity of $\phi_\ell(\cdot)$ and the condition $\mathbf{A}_\ell y + b_\ell \in K_\ell^*$ whenever $y \in Y$. The associated Nash v.i. is given by the domain Z and the monotone mapping

$$(4.9) \quad F(z) \equiv F(x, y) = \left[\sum_{\ell=1}^m [\phi'_\ell(x)]^* [\mathbf{A}_\ell y + b_\ell]; - \sum_{\ell=1}^m \mathbf{A}_\ell^* \phi_\ell(x) + \Phi'_*(y) \right].$$

The advantage of the v.i. reformulation of (4.1) is that F is linear in $\phi_\ell(\cdot)$, so that the initial unbiased SO for ϕ_ℓ induces an unbiased stochastic oracle for F , specifically, the oracle

$$(4.10) \quad \Xi(x, y, \zeta_i) = \left[\sum_{\ell=1}^m \mathbf{G}_\ell^*(x, \zeta_i) [\mathbf{A}_\ell y + b_\ell]; - \sum_{\ell=1}^m \mathbf{A}_\ell^* f_\ell(x, \zeta_i) + \Phi'_*(y) \right].$$

We are about to use this oracle in order to solve the stochastic composite minimization problem (4.1) by the SMP algorithm.

4.2. Setup for the SMP as applied to (4.9). In retrospect, the setup for SMP we are about to present is a kind of the best – resulting in the best possible efficiency estimate (3.12) – we can build from the entities participating in the description of the problem (4.1). Specifically, we equip the space $\mathcal{E} = \mathcal{X} \times \mathcal{Y}$ with the norm

$$\|[(x, y)]\| \equiv \sqrt{\|x\|_x^2 / \Omega_x^2 + \|y\|_y^2 / \Omega_y^2},$$

the conjugate norm clearly is

$$\|(\xi, \eta)\|_* = \sqrt{\Omega_x^2 \|\xi\|_{x,*}^2 + \Omega_y^2 \|\eta\|_{y,*}^2}.$$

Finally, we equip $Z = X \times Y$ with the distance-generating function

$$\omega(x, y) = \frac{1}{\alpha_x \Omega_x^2} \omega_x(x) + \frac{1}{\alpha_y \Omega_y^2} \omega_y(y).$$

The SMP-related properties of our setup are summarized in the following

LEMMA 4.1. *Let*

$$(4.11) \quad \mathcal{A} = \max_{y \in \mathcal{Y}: \|y\|_y \leq 1} \sum_{\ell=1}^m \|\mathbf{A}_\ell y\|_{(\ell,*)}, \quad \mathcal{B} = \sum_{\ell=1}^m \|b_\ell\|_{(\ell,*)}.$$

(i) *The parameters of the just defined distance-generating function ω w.r.t. the just defined norm $\|\cdot\|$ are $\alpha = 1$, $\Theta = 1$, $\Omega = \sqrt{2}$.*

(ii) *One has*

$$(4.12) \quad \forall(z, z' \in Z) : \|F(z) - F(z')\|_* \leq L \|z - z'\| + M,$$

where

$$\begin{aligned} L &= 5\mathcal{A}\Omega_x\Omega_y[\Omega_x L_x + M_x] + \mathcal{B}\Omega_x^2 L_x + \Omega_y^2 L_y \\ M &= [2\mathcal{A}\Omega_y + \|b\|_1]\Omega_x M_x + \Omega_y M_y \end{aligned}$$

Besides this,

$$(4.13) \quad \forall(z \in Z, i) : \mathbf{E}\{\Xi(z, \zeta_i)\} = F(z); \quad \mathbf{E}\{\|\Xi(z, \zeta_i) - F(z)\|_*^2\} \leq M^2.$$

Furthermore, if relations (4.3.b,d) are strengthened to

$$(4.14) \quad \mathbf{E}\left\{\exp\left\{\max_{1 \leq \ell \leq m} \|f_\ell(x, \zeta_i) - \phi_\ell(x)\|_{(\ell)}^2 / (\Omega_x M)^2\right\}\right\} \leq \exp\{1\},$$

$$\mathbf{E}\left\{\exp\left\{\max_{\substack{h \in \mathcal{X}, \\ \|h\|_x \leq 1}} \|[\mathbf{G}_\ell(x) - \phi'_\ell(x)]h\|_{(\ell)}^2 / M^2\right\}\right\} \leq \exp\{1\}, \quad 1 \leq \ell \leq m,$$

then

$$(4.15) \quad \mathbf{E}\{\exp\{\|\Xi(z, \zeta_i) - F(z)\|_*^2 / M^2\}\} \leq \exp\{1\}.$$

Combining Lemma 4.1 with Corollary (3.6) we get explicit efficiency estimates for the SMP algorithm as applied to the Stochastic composite minimization problem (4.1).

4.3. Application to Stochastic Semidefinite Feasibility problem. Assume we are interested to solve a feasible system of matrix inequalities

$$(4.16) \quad \psi_\ell(x) \preceq 0, \quad \ell = 1, \dots, m \ \& \ x \in X,$$

where $m > 1$, $X \subset \mathcal{X}$ is as in the description of the Stochastic composite problem, and $\psi_\ell(\cdot)$ take values in the spaces $\mathcal{E}_\ell = \mathbf{S}^{p_\ell}$ of symmetric $p_\ell \times p_\ell$ matrices. We equip \mathcal{E}_ℓ with the Frobenius inner product, the semidefinite cone $K_\ell = \mathbf{S}_+^{p_\ell}$ and the spectral norm $\|\cdot\|_{(\ell)} = |\cdot|_\infty$ (recall that $|A|_\infty$ is the maximal singular value of matrix A). We assume that ψ_ℓ are Lipschitz continuous and $K_\ell = \mathbf{S}_+^{p_\ell}$ -convex functions on X such that for all $x, x' \in X$ and for all ℓ one has

$$(4.17) \quad \begin{aligned} \max_{h \in \mathcal{X}, \|h\|_x \leq 1} |[\psi'_\ell(x) - \psi'_\ell(x')]h|_\infty &\leq L_\ell \|x - x'\|_{(\ell)} + M_\ell, \\ \max_{h \in \mathcal{X}, \|h\|_x \leq 1} |\psi'_\ell(x)h|_\infty &\leq L_\ell \Omega_x + M_\ell \end{aligned}$$

for certain selections $\psi'_\ell(x) \in \partial^{K_\ell} \psi_\ell(x)$, $x \in X$, with some known nonnegative constants L_ℓ, M_ℓ .

We assume that $\psi_\ell(\cdot)$ are represented by an SO which at i -th call, the input being $x \in X$, returns the matrices $\widehat{f}_\ell(x, \zeta_i) \in \mathbf{S}^{p_\ell}$ and the linear maps $\widehat{\mathbf{G}}_\ell(x, \zeta_i)$ from \mathcal{X} to \mathcal{E}_ℓ such that for all $x \in X$ it holds

$$(4.18) \quad \begin{aligned} (a) \quad &\mathbf{E}\left\{\widehat{f}_\ell(x, \zeta_i)\right\} = \psi_\ell(x), \quad \mathbf{E}\left\{\widehat{\mathbf{G}}_\ell(x, \zeta_i)\right\} = \psi'_\ell(x), \quad 1 \leq \ell \leq m \\ (b) \quad &\mathbf{E}\left\{\max_{1 \leq \ell \leq m} |\widehat{f}_\ell(x, \zeta_i) - \psi_\ell(x)|_\infty^2 / (\Omega_x M_\ell)^2\right\} \leq 1 \\ (c) \quad &\mathbf{E}\left\{\max_{\substack{h \in \mathcal{X}, \\ \|h\|_x \leq 1}} |[\widehat{\mathbf{G}}_\ell(x, \zeta_i) - \psi'_\ell(x)]h|_\infty^2 / M_\ell^2\right\} \leq 1, \quad 1 \leq \ell \leq m. \end{aligned}$$

Given a number t of steps of the SMP algorithm, let us act as follows.

A. We compute the m quantities $\mu_\ell = \frac{\Omega_x L_\ell}{\sqrt{t}} + M_\ell$, $\ell = 1, \dots, m$, and set

$$(4.19) \quad \mu = \max_{1 \leq \ell \leq m} \mu_\ell, \quad \beta_\ell = \frac{\mu}{\mu_\ell}, \quad \phi_\ell(\cdot) = \beta_\ell \psi_\ell(\cdot), \quad L_x = \frac{\mu \sqrt{t}}{\Omega_x}, \quad M_x = \mu.$$

Note that by construction $\beta_\ell \geq 1$ and $L_x/L_\ell \geq \beta_\ell$, $M_x/M_\ell \geq \beta_\ell$ for all ℓ , so that the functions ϕ_ℓ satisfy (4.2) with the just defined L_x , M_x . Further, the SO for $\psi_\ell(\cdot)$'s can be converted into an SO for $\phi_\ell(\cdot)$'s by setting

$$f_\ell(x, \zeta) = \beta_\ell \widehat{f}_\ell(x, \zeta), \quad \mathbf{G}_\ell(x, \zeta) = \beta_\ell \widehat{\mathbf{G}}_\ell(x, \zeta).$$

By (4.18), this oracle satisfies (4.3).

B. We then build the Stochastic Matrix Minimax problem

$$(4.20) \quad \min_{x \in X} \max_{1 \leq \ell \leq m} \lambda_{\max}(\phi_\ell(x)),$$

associated with the just defined ϕ_1, \dots, ϕ_m , that is, the Stochastic composite problem (4.1) associated with ϕ_1, \dots, ϕ_m and the outer function

$$\begin{aligned} \Phi(u_1, \dots, u_m) &= \max_{1 \leq \ell \leq m} \lambda_{\max}(u_\ell) = \max_{y \in Y} \sum_{\ell=1}^m \langle u_\ell, y_\ell \rangle_F, \\ Y &= \{y = \text{Diag}\{y_1, \dots, y_m\} \in \mathcal{Y} = \mathbf{S}^{p_1} \times \dots \times \mathbf{S}^{p_m} : y \succeq 0, \text{Tr}(y) = 1\} \\ &\subset \mathcal{Y} = \mathbf{S}^{p_1} \times \dots \times \mathbf{S}^{p_m}. \end{aligned}$$

Thus in the notation from (4.4) we have $\mathbf{A}_\ell y = y_\ell$, $b_\ell = 0$, $\Phi_* \equiv 0$. Hence $L_x = M_x = 0$, and Y is a spectahedron. We equip \mathcal{Y} and Y with the Spectahedron setup, arriving at

$$\alpha_y = 1, \quad \Theta_y = \ln \sum_{\ell=1}^m p_\ell, \quad \Omega_y = \sqrt{2 \ln \sum_{\ell=1}^m p_\ell}.$$

C. We have specified all entities participating in the description of the Stochastic composite problem. It is immediately seen that these entities satisfy all conditions of Section 4.1. We can now solve the resulting Stochastic composite problem by t -step SMP algorithm with the setup presented in Section 4.2. The corresponding convex-concave saddle point problem is

$$\min_{x \in X} \max_{y \in Y} \sum_{\ell=1}^m \beta_\ell \langle \psi_\ell(x), y_\ell \rangle_F;$$

with the monotone operator and SO, respectively,

$$\begin{aligned} F(z) &\equiv F(x, y) = \left[\sum_{\ell=1}^m \beta_\ell [\psi'_\ell(x)]^* y_\ell; - \text{Diag} \{ \alpha_1 \psi_1(x), \dots, \alpha_m \psi_m(x) \} \right], \\ \Xi((x, y), \zeta) &= \left[\sum_{\ell=1}^m \beta_\ell \widehat{\mathbf{G}}_\ell^*(x, \zeta) y_\ell; - \text{Diag} \{ \alpha_1 \widehat{f}_1(x, \zeta), \dots, \alpha_m \widehat{f}_m(x, \zeta) \} \right]. \end{aligned}$$

Combining Lemma 4.1, Corollary 3.6 and taking into account the origin of the quantities L_x , M_x , and that $\mathcal{A} = 1$, $\mathcal{B} = 0$ ³⁾, we arrive at the following result:

³⁾ See (4.11) and note that we are in the case when $b_\ell = 0$ and $\|\cdot\|_{(\ell,*)}$ is the trace norm; thus, $\sum_{\ell=1}^m \|\mathbf{A}_\ell y\|_{(\ell,*)} = \sum_{\ell=1}^m |y_\ell|_1 = |y|_1 = \|y\|_y$.

PROPOSITION 4.2. *With the outlined construction, the resulting s.v.i. reads*

$$(4.21) \quad \text{find } z_* \in Z = X \times Y : \langle F(z), z - z_* \rangle \geq 0 \quad \forall z \in Z,$$

for the monotone operator F which satisfies (1.4) with

$$L = 10 \left[\ln \sum_{\ell=1}^m p_\ell \right]^{\frac{1}{2}} \Omega_x \mu (\sqrt{t} + 1), \quad M = 4 \left[\ln \sum_{\ell=1}^m p_\ell \right]^{\frac{1}{2}} \Omega_x \mu;$$

Beside this, the resulting SO for F satisfies (4.13) with the just defined value of M .

Let now

$$\gamma = \left(10 \left[3 \ln \sum_{\ell=1}^m p_\ell \right]^{\frac{1}{2}} \Omega_x \mu (\sqrt{t} + 1) \right)^{-1}, \quad 1 \leq \tau \leq t.$$

When applying to (4.21) the t -step SMP algorithm with the constant stepsizes $\gamma_\tau \equiv \gamma$ (cf. (3.11) and note that we are in the situation $\alpha = \Theta = 1$), we get an approximate solution $\hat{z}_t = (\hat{x}_t, \hat{y}_t)$ such that

$$(4.22) \quad \mathbf{E} \left\{ \max_{1 \leq \ell \leq m} \beta_\ell \lambda_{\max}(\psi_\ell(\hat{x}_t)) \right\} \leq 80 \frac{\Omega_x [\ln \sum_{\ell=1}^m p_\ell]^{\frac{1}{2}} \mu}{\sqrt{t}}$$

(cf. (3.12) and take into account that we are in the case of $\Omega = \sqrt{2}$, while the optimal value in (4.20) is nonpositive, since (4.16) is feasible).

Furthermore, if assumptions (4.18.b,c) are strengthened to

$$\begin{aligned} \mathbf{E} \left\{ \max_{1 \leq \ell \leq m} \exp\{|\hat{f}_\ell(x, \zeta_i) - \psi_\ell(x)|_\infty^2 / (\Omega_x M_\ell)^2\} \right\} &\leq \exp\{1\}, \\ \mathbf{E} \left\{ \exp\left\{ \max_{h \in \mathcal{X}, \|h\|_x \leq 1} |[\hat{\mathbf{G}}_\ell(x, \zeta_i) - \psi'_\ell(x)]h|_\infty^2 / M_\ell^2 \right\} \right\} &\leq \exp\{1\}, \quad 1 \leq \ell \leq m, \end{aligned}$$

then, in addition to (4.22), we have for any $\Lambda > 0$:

$$\begin{aligned} \text{Prob} \left\{ \max_{1 \leq \ell \leq m} \beta_\ell \lambda_{\max}(\psi_\ell(\hat{x}_t)) > 80 \frac{\Omega_x [\ln \sum_{\ell=1}^m p_\ell]^{\frac{1}{2}} \mu}{\sqrt{t}} + \Lambda \frac{15 [\ln \sum_{\ell=1}^m p_\ell]^{\frac{1}{2}} \mu}{\sqrt{t}} \right\} \\ \leq \exp\{-\Lambda^2/3\} + \exp\{-\Lambda t\}. \end{aligned}$$

Discussion. Imagine that instead of solving the system of matrix inequalities (4.16), we were interested to solve just a single matrix inequality $\psi_\ell(x) \preceq 0$, $x \in X$. When solving this inequality by the SMP algorithm as explained above, the efficiency estimate would be

$$\begin{aligned} \mathbf{E} \{ \psi_\ell(\hat{x}_t^\ell) \} &\leq O(1) [\ln(p_\ell + 1)]^{1/2} \Omega_x \left[\frac{\Omega_x L_\ell}{t} + \frac{M_\ell}{\sqrt{t}} \right] = O(1) [\ln(p_\ell + 1)]^{1/2} \Omega_x \frac{\mu_\ell}{\sqrt{t}} \\ &= O(1) [\ln(p_\ell + 1)]^{1/2} \beta_\ell^{-1} \frac{\Omega_x \mu}{\sqrt{t}}, \end{aligned}$$

(recall that the matrix inequality in question is feasible), where \hat{x}_t^ℓ is the resulting approximate solution. Looking at (4.22), we see that the expected accuracy of the SMP as applied, in the aforementioned manner, to (4.16) is only by a logarithmic in $\sum_\ell p_\ell$ factor worse:

$$(4.23) \quad \mathbf{E} \{\psi_\ell(\hat{x}_t)\} \leq O(1) \left[\ln \sum_{\ell=1}^m p_\ell \right]^{1/2} \beta_\ell^{-1} \frac{\Omega_x \mu_\ell}{\sqrt{t}} = O(1) \left[\ln \sum_{\ell=1}^m p_\ell \right]^{1/2} \frac{\Omega_x \mu_\ell}{\sqrt{t}}.$$

Thus, as far as the quality of the SPM-generated solution is concerned, passing from solving a single matrix inequality to solving a system of m inequalities is “nearly costless”. As an illustration, consider the case where some of ψ_ℓ are “easy” – smooth and easy-to-observe ($M_\ell = 0$), while the remaining ψ_ℓ are “difficult”, i.e., might be non-smooth and/or difficult-to-observe ($L_\ell = 0$). In this case, (4.23) reads

$$\mathbf{E} \{\psi_\ell(\hat{x}_t)\} \leq \begin{cases} O(1) \left[\ln \sum_{\ell=1}^m p_\ell \right]^{1/2} \frac{\Omega_x^2 L_\ell}{t}, & \psi_\ell \text{ is easy,} \\ O(1) \left[\ln \sum_{\ell=1}^m p_\ell \right]^{1/2} \frac{\Omega_x M_\ell}{\sqrt{t}}, & \psi_\ell \text{ is difficult.} \end{cases}$$

In other words, the violations of the easy and the difficult constraints in (4.16) converge to 0 as $t \rightarrow \infty$ with the rates $O(1/t)$ and $O(1/\sqrt{t})$, respectively. It should be added that when X is the unit Euclidean ball in $\mathcal{X} = \mathbf{R}^n$ and X, \mathcal{X} are equipped with the Euclidean setup, the rates of convergence $O(1/t)$ and $O(1/\sqrt{t})$ are the best rates one can achieve without imposing bounds on n and/or imposing additional restrictions on ψ_ℓ 's.

4.4. Eigenvalue optimization via SMP. The problem we are interested in now is

$$(4.24) \quad \begin{aligned} \text{Opt} &= \min_{x \in X} f(x) := \lambda_{\max}(A_0 + x_1 A_1 + \dots + x_n A_n), \\ X &= \{x \in \mathbf{R}^n : x \geq 0, \sum_{i=1}^n x_i = 1\}, \end{aligned}$$

where A_0, A_1, \dots, A_n , $n > 1$, are given symmetric matrices with common block-diagonal structure (p_1, \dots, p_m) . I.e., all A_j are block-diagonal with diagonal blocks A_j^ℓ of sizes $p_\ell \times p_\ell$, $1 \leq \ell \leq m$. We denote

$$p^{(\kappa)} = \sum_{\ell=1}^m p_\ell^\kappa, \quad \kappa = 1, 2, 3; \quad p^{\max} = \max_\ell p_\ell.$$

Setting

$$\phi_\ell : X \mapsto \mathcal{E}_\ell = \mathbf{S}^{p_\ell}, \quad \phi_\ell(x) = A_0^\ell + \sum_{j=1}^n x_j A_j^\ell, \quad 1 \leq \ell \leq m,$$

we represent (4.24) as a particular case of the Matrix Minimax problem (4.6), with all functions $\phi_\ell(x)$ being affine and X being the standard simplex in $\mathcal{X} = \mathbf{R}^n$.

Now, since A_j are known in advance, there is nothing stochastic in our problem, and it can be solved either by interior point methods, or by “computationally cheap” gradient-type methods; these latter methods are preferable when the problem is large-scale and medium accuracy solutions are sought. For instance, one can apply the t -step

(deterministic) Mirror Prox algorithm from [6] to the saddle point reformulation (4.8) of our specific Matrix Minimax problem, i.e., to the saddle point problem

$$(4.25) \quad \min_{x \in X} \max_{y \in Y} \langle y, A_0 + \sum_{j=1}^n x_j A_j \rangle_F, \\ Y = \{y = \text{Diag}\{y_1, \dots, y_m\} : y_\ell \in \mathbf{S}_+^{p_\ell}, 1 \leq \ell \leq m, \text{Tr}(Y) = 1\}.$$

The accuracy of the approximate solution \tilde{x}_t of the (deterministic) Mirror Prox algorithm is [6, Example 2]

$$f(\tilde{x}_t) - \text{Opt} \leq O(1) \frac{\sqrt{\ln(n) \ln(p^{(1)})} A_\infty}{t}.$$

This efficiency estimate is the best known so far among those attainable with “computationally cheap” deterministic methods. On the other hand, the complexity of one step of the algorithm is dominated, up to an absolute constant factor, by the necessity, given $x \in X$ and $y \in Y$,

1. to compute the matrix $A_0 + \sum_{j=1}^n x_j A_j$ and the vector $[\text{Tr}(Y A_1); \dots; \text{Tr}(Y A_n)]$;
2. to compute the eigenvalue decomposition of y .

When using the standard Linear Algebra, the computational effort per step is

$$C_{det} = O(1)[np^{(2)} + p^{(3)}]$$

arithmetic operations.

We are about to demonstrate that one can equip the deterministic problem in question by an “artificial” SO in such a way that the associated SMP algorithm, under certain circumstances, exhibits better performance than deterministic algorithms. Let us consider the following construction of the SO for F (different from the SO (4.10)!). Observe that the monotone operator associated with the saddle point problem (4.25) is

$$(4.26) \quad F(x, y) = \left[\underbrace{[\text{Tr}(y A_1); \dots; \text{Tr}(y A_n)]}_{F^x(x, y)}; \underbrace{-A_0 - \sum_{j=1}^n x_j A_j}_{F^y(x, y)} \right].$$

Given $x \in X$, $y = \text{Diag}\{y_1, \dots, y_m\} \in Y$, we build a random estimate $\Xi = [\Xi^x; \Xi^y]$ of $F(x, y) = [F^x(x, y); F^y(x, y)]$ as follows:

1. we generate a realization j of a random variable taking values $1, \dots, n$ with probabilities x_1, \dots, x_n (recall that $x \in X$, the standard simplex, so that x indeed can be seen as a probability distribution), and set

$$(4.27) \quad \Xi^y = A_0 + A_j;$$

2. we compute the quantities $\nu_\ell = \text{Tr}(y_\ell)$, $1 \leq \ell \leq m$. Since $y \in Y$, we have $\nu_\ell \geq 0$ and $\sum_{\ell=1}^m \nu_\ell = 1$. We further generate a realization i of random variable taking values $1, \dots, m$ with probabilities ν_1, \dots, ν_m , and set

$$(4.28) \quad \Xi^x = [\text{Tr}(A_1^i \bar{y}_i); \dots; \text{Tr}(A_n^i \bar{y}_i)], \quad \bar{y}_i = (\text{Tr}(y_i))^{-1} y_i.$$

The just defined random estimate Ξ of $F(x, y)$ can be expressed as a deterministic function $\Xi(x, y, \eta)$ of (x, y) and random variable η uniformly distributed on $[0, 1]$. Given x, y and η , the value of this function can be computed with the arithmetic cost

$O(1)(n(p^{\max})^2 + p^{(2)})$ (indeed, $O(1)(n + p^{(1)})$ operations are needed to convert η into i and j , $O(1)p^{(2)}$ operations are used to write down the y -component $-A_0 - A_j$ of Ξ , and $O(1)n(p^{\max})^2$ operations are needed to compute Ξ^x). Now consider the SO's Ξ_k (k is a positive integer) obtained by averaging the outputs of k calls to our basic oracle Ξ . Specifically, at the i -t call to the oracle Ξ_k , $z = (x, y) \in Z = X \times Y$ being the input, the oracle returns the vector

$$\Xi_k(z, \zeta_i) = \frac{1}{k} \sum_{s=1}^k \Xi(z, \eta_{is}),$$

where $\zeta_i = [\eta_{i1}; \dots; \eta_{ik}]$ and $\{\eta_{is}\}_{1 \leq i, 1 \leq s \leq k}$ are independent random variables uniformly distributed on $[0, 1]$. Note that the arithmetic cost of a single call to Ξ_k is

$$\mathcal{C}_k = O(1)k(n(p^{\max})^2 + p^{(2)}).$$

The Nash v.i. associated with the saddle point problem (4.25) with the stochastic oracle Ξ_k (k being the first parameter of our construction) specify a Nash s.v.i. on the domain $Z = X \times Y$. Let us equip the standard simplex X and its embedding space $\mathcal{X} = \mathbf{R}^n$ with the Simplex setup, and the spectahedron Y and its embedding space $\mathcal{Y} = \mathbf{S}^{p_1} \times \dots \times \mathbf{S}^{p_m}$ with the Spectahedron setup (see Section 2.2). Let us next combine the x - and the y -setups, exactly as explained in the beginning of Section 4.2, into an SMP setup for the domain $Z = X \times Y$ – a distance-generating function $\omega(\cdot)$ and a norm $\|\cdot\|$ on the embedding space $\mathbf{R}^n \times (\mathbf{S}^{p_1} \times \dots \times \mathbf{S}^{p_\ell})$ of Z . The SMP-related properties of the resulting setup are summarized in the following statement.

LEMMA 4.3. *Let $n \geq 3$, $p^{(1)} \geq 3$. Then*

(i) *The parameters of the just defined distance-generating function ω w.r.t. the just defined norm $\|\cdot\|$ are $\alpha = 1$, $\Theta = 1$, $\Omega = \sqrt{2}$.*

(ii) *For any $z, z' \in Z$ one has*

$$(4.29) \quad \|F(z) - F(z')\|_* \leq L\|z - z'\|, \quad L = 2 \ln(n) + 4 \ln(p^{(1)}).$$

Besides this, for any ($z \in Z$, $i = 1, 2, \dots$),

$$(4.30) \quad \begin{aligned} (a) \quad & \mathbf{E} \{ \Xi_k(z, \zeta_i) \} = F(z); \\ (b) \quad & \mathbf{E} \left\{ \exp \left\{ \frac{\| \Xi(z, \zeta_i) - F(z) \|_*^2}{M^2} \right\} \right\} \leq \exp\{1\}, \\ & M = 27[\ln(n) + \ln(p^{(1)})]A_\infty / \sqrt{k}. \end{aligned}$$

REFERENCES

- [1] Azuma, K. Weighted sums of certain dependent random variables. *Tökoku Math. J.*, **19** (1967), 357-367.
- [2] Ben-Tal, A., Nemirovski, A. "Non-Euclidean restricted memory level method for large-scale convex optimization" – *Math. Progr.* **102** (2005), 407-456.
- [3] Juditsky, A. Lan, G., Nemirovski, A., Shapiro, A., Stochastic Approximation Approach to Stochastic Programming, http://www.optimization-online.org/DB_HTML/2007/09/1787.html
- [4] Juditsky, A., Nemirovski, A. (2008), Large Deviations of Vector-valued Martingales in 2-Smooth Normed Spaces
E-print: http://www.optimization-online.org/DB_HTML/2008/04/1947.html
- [5] Nemirovski, A., Yudin, D., *Problem complexity and method efficiency in Optimization* J. Wiley & Sons (1983).

- [6] A. Nemirovski, “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** (2004), 229-251.
- [7] Lu, Z., Nemirovski A., Monteiro, R. “Large-Scale Semidefinite Programming via Saddle Point Mirror-Prox Algorithm”, *Math. Progr.*, **109** (2007), 211-237.
- [8] Nemirovski, A., Onn, S., Rothblum, U. (2007), “Accuracy certificates for computational problems with convex structure” – submitted to *Mathematics of Operations Research*. E-print: http://www.optimization-online.org/DB_HTML/2007/04/1634.html
- [9] Nesterov, Yu. “Smooth minimization of non-smooth functions”, *Math. Progr.*, **103** (2005), 127-152.
- [10] Nesterov, Yu. “Excessive gap technique in nonsmooth convex minimization”, *SIAM J. Optim.*, **16** (2005), 235-249.
- [11] Nesterov, Yu. “Dual extrapolation and its applications to solving variational inequalities and related problems”, *Math. Progr.* **109** (2007) 319-344.

5. Appendix.

5.1. Proof of Theorem 3.2. We start with the following simple observation: if r_e is a solution to (2.2), then $\partial_Z \omega(r_e)$ contains $-e$ and thus is nonempty, so that $r_e \in Z^\circ$. Moreover, one has

$$(5.1) \quad \langle \omega'(r_e) - e, u - r_e \rangle \geq 0 \quad \forall u \in Z.$$

Indeed, by continuity argument, it suffices to verify the inequality in the case when $u \in \text{rint}(Z) \subset Z^\circ$. For such an u , the convex function

$$f(t) = \omega(r_e + t(u - r_e)) + \langle r_e + t(u - r_e), e \rangle, \quad t \in [0, 1]$$

is continuous on $[0, 1]$ and has a continuous on $[0, 1]$ field of subgradients

$$g(t) = \langle \omega'(r_e + t(u - r_e)) + e, u - r_e \rangle.$$

It follows that the function is continuously differentiable on $[0, 1]$ with the derivative $g(t)$. Since the function attains its minimum on $[0, 1]$ at $t = 0$, we have $g(0) \geq 0$, which is exactly (5.1).

At least the first statement of the following Lemma is well-known:

LEMMA 5.1. *For every $z \in Z^\circ$, the mapping $\xi \mapsto P(z, \xi)$ is a single-valued mapping of \mathcal{E} onto Z° , and this mapping is Lipschitz continuous, specifically,*

$$(5.2) \quad \|P(z, \zeta) - P(z, \eta)\| \leq \alpha^{-1} \|\zeta - \eta\|_* \quad \forall \zeta, \eta \in \mathcal{E}.$$

Besides this,

$$(5.3) \quad \begin{aligned} (a) \quad \forall (u \in Z) : V(P(z, \zeta), u) &\leq V(z, u) + \langle \zeta, u - P(z, \zeta) \rangle - V(z, P(z, \zeta)) \\ (b) &\leq V(z, u) + \langle \zeta, u - z \rangle + \frac{\|\zeta\|_*^2}{2\alpha}. \end{aligned}$$

Proof.

Let $v \in P(z, \zeta)$, $w \in P(z, \eta)$. As $V'_u(z, u) = \omega'(u) - \omega'(z)$, invoking 5.1, we have $v, w \in Z^\circ$ and

$$(5.4) \quad \langle \omega'(v) - \omega'(z) + \zeta, v - u \rangle \leq 0 \quad \forall u \in Z.$$

$$(5.5) \quad \langle \omega'(w) - \omega'(z) + \eta, w - u \rangle \leq 0 \quad \forall u \in Z.$$

Setting $u = w$ in (5.4) and $u = v$ in (5.5), we get

$$\langle \omega'(v) - \omega'(z) + \zeta, v - w \rangle \leq 0, \quad \langle \omega'(w) - \omega'(z) + \eta, v - w \rangle \geq 0,$$

whence $\langle \omega'(w) - \omega'(v) + [\eta - \zeta], v - w \rangle \geq 0$, or

$$\|\eta - \zeta\|_* \|v - w\| \geq \langle \eta - \zeta, v - w \rangle \geq \langle \omega'(v) - \omega'(w), v - w \rangle \geq \alpha \|v - w\|^2,$$

and (5.2) follows. This relation, as a byproduct, implies that $P(z, \cdot)$ is single-valued.

To prove (5.3), let $v = P(z, \zeta)$. We have

$$\begin{aligned} V(v, u) - V(z, u) &= [\omega(u) - \langle \omega'(v), u - v \rangle - \omega(v)] - [\omega(u) - \langle \omega'(z), u - z \rangle - \omega(z)] \\ &= \langle \omega'(v) - \omega'(z) + \zeta, v - u \rangle + \langle \zeta, u - v \rangle - [\omega(v) - \langle \omega'(z), v - z \rangle - \omega(z)] \\ (\text{due to (5.4)}) &\leq \langle \zeta, u - v \rangle - V(z, v), \end{aligned}$$

as required in (a) of (5.3). The bound (b) of (5.3) is obtained from (5.3) using the Young inequality:

$$\langle \zeta, z - v \rangle \leq \frac{\|z\|_*^2}{2\alpha} + \frac{\alpha}{2} \|z - v\|^2.$$

Indeed, observe that by definition, $V(z, \cdot)$ is strongly convex with parameter α , and $V(z, v) \geq \frac{\alpha}{2} \|z - v\|^2$, so that

$$\langle \zeta, u - v \rangle - V(z, v) = \langle \zeta, u - z \rangle + \langle \zeta, z - v \rangle - V(z, v) \leq \langle \zeta, u - z \rangle + \frac{\|\zeta\|_*^2}{2\alpha}.$$

□

We have the following simple corollary of Lemma 5.1:

COROLLARY 5.2. *Let ξ_1, ξ_2, \dots be a sequence of elements of \mathcal{E} . Define the sequence $\{y_\tau\}_{\tau=0}^\infty$ in Z° as follows:*

$$y_\tau = P(y_{\tau-1}, \xi_\tau), \quad y_0 \in Z^\circ.$$

Then y_τ is a measurable function of y_0 and ξ_1, \dots, ξ_τ such that

$$(5.6) \quad (\forall u \in Z) : \quad \left\langle -\sum_{\tau=1}^t \xi_\tau, u \right\rangle \leq V(y_0, u) + \sum_{\tau=1}^t \zeta_\tau,$$

with

$$(5.7) \quad |\zeta_\tau| \leq r \|\xi_\tau\|_* \quad (\text{here } r = \max_{u \in Z} \|u\|); \quad \zeta_\tau \leq -\langle \xi_\tau, y_{\tau-1} \rangle + \frac{\|\xi_\tau\|_*^2}{2\alpha}.$$

Proof. Using the bound (b) of (5.3) with $\zeta = \xi_t$ and $z = y_{t-1}$ (so that $y_t = P(y_{t-1}, \xi_t)$) we obtain for any $u \in Z$:

$$V(y_t, u) - V(y_{t-1}, u) - \langle \xi_t, u \rangle \leq -\langle \xi_t, y_t \rangle - V(y_{t-1}, y_t) \equiv \zeta_t.$$

Note that

$$\zeta_t = \max_{v \in Z} [-\langle \xi_t, v \rangle - V(y_{t-1}, v)],$$

so that

$$-r \|\xi_t\|_* \leq -\langle \xi_t, y_{t-1} \rangle \leq \zeta_t \leq r \|\xi_t\|_*.$$

Further, due to the strong convexity of V ,

$$\zeta_t = -\langle \xi_t, y_{t-1} \rangle + [-\langle \xi_t, y_t - y_{t-1} \rangle - V(y_{t-1}, y_t)] \leq -\langle \xi_t, y_{t-1} \rangle + \frac{\|\xi_t\|_*^2}{2\alpha}.$$

When summing up from $\tau = 1$ to $\tau = t$ we arrive at the corollary. \square

We also need the following result.

LEMMA 5.3. *Let $z \in Z^o$, let ζ, η be two points from \mathcal{E} , and let*

$$w = P(z, \zeta), \quad r_+ = P(z, \eta)$$

Then for all $u \in Z$ one has

$$(5.8) \quad \begin{aligned} (a) \quad & \|w - r_+\| \leq \alpha^{-1} \|\zeta - \eta\|_* \\ (b) \quad & V(r_+, u) - V(z, u) \leq \langle \eta, u - w \rangle + \frac{\|\zeta - \eta\|_*^2}{2\alpha} - \frac{\alpha}{2} \|w - z\|^2. \end{aligned}$$

Proof. (a): this is nothing but (5.2).

(b): Using (a) of (5.3) in Lemma 5.1 we can write for $u = r_+$:

$$V(w, r_+) \leq V(z, r_+) + \langle \zeta, r_+ - w \rangle - V(z, w).$$

This results in

$$(5.9) \quad V(z, r_+) \geq V(w, r_+) + V(z, w) + \langle \zeta, w - r_+ \rangle.$$

Using (5.3) with η substituted for ζ we get

$$\begin{aligned} V(r_+, u) & \leq V(z, u) + \langle \eta, u - r_+ \rangle - V(z, r_+) \\ & = V(z, u) + \langle \eta, u - w \rangle + \langle \eta, w - r_+ \rangle - V(z, r_+) \\ \text{[by (5.9)]} & \leq V(z, u) + \langle \eta, u - w \rangle + \langle \eta - \zeta, w - r_+ \rangle - V(z, w) - V(w, r_+) \\ & \leq V(z, u) + \langle \eta, u - w \rangle + \langle \eta - \zeta, w - r_+ \rangle - \frac{\alpha}{2} [\|w - z\|^2 + \|w - r_+\|^2], \end{aligned}$$

due to the strong convexity of V . To conclude the bound (b) of (5.8) it suffices to note that by the Young inequality,

$$\langle \eta - \zeta, w - r_+ \rangle \leq \frac{\|\eta - \zeta\|_*^2}{2\alpha} + \frac{\alpha}{2} \|w - r_+\|^2.$$

\square

We are able now to prove Theorem 3.2. By (1.4) we have that

$$(5.10) \quad \begin{aligned} \|\widehat{F}(w_\tau) - \widehat{F}(r_{\tau-1})\|_*^2 & \leq (L\|r_{\tau-1} - w_\tau\| + M + \epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2 \\ & \leq 3L^2\|w_\tau - r_{\tau-1}\|^2 + 3M^2 + 3(\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2. \end{aligned}$$

Let us now apply Lemma 5.3 with $z = r_{\tau-1}$, $\zeta = \gamma_\tau \widehat{F}(r_{\tau-1})$, $\eta = \gamma_\tau \widehat{F}(w_\tau)$ (so that $w = w_\tau$ and $r_+ = r_\tau$). We have for any $u \in Z$

$$\begin{aligned} & \langle \gamma_\tau \widehat{F}(w_\tau), w_\tau - u \rangle + V(r_\tau, u) - V(r_{\tau-1}, u) \\ & \leq \frac{\gamma_\tau^2}{2\alpha} \|\widehat{F}(w_\tau) - \widehat{F}(r_{\tau-1})\|_*^2 - \frac{\alpha}{2} \|w_\tau - r_{\tau-1}\|^2 \\ \text{[by (5.10)]} & \leq \frac{3\gamma_\tau^2 L^2}{2\alpha} [\|w_\tau - r_{\tau-1}\|^2 + M^2 + (\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2] - \frac{\alpha}{2} \|w_\tau - r_{\tau-1}\|^2 \\ \text{[by (3.4)]} & \leq \frac{3\gamma_\tau^2}{2\alpha} [M^2 + (\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2] \end{aligned}$$

When summing up from $\tau = 1$ to $\tau = t$ we obtain

$$\begin{aligned} \sum_{\tau=1}^t \langle \gamma_\tau \widehat{F}(w_\tau), w_\tau - u \rangle &\leq V(r_0, u) - V(r_t, u) + \sum_{\tau=1}^t \frac{3\gamma_\tau^2}{2\alpha} [M^2 + (\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2] \\ &\leq \Theta(r_0) + \sum_{\tau=1}^t \frac{3\gamma_\tau^2}{2\alpha} [M^2 + (\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2]. \end{aligned}$$

Hence, for all $u \in Z$,

$$\begin{aligned} &\sum_{\tau=1}^t \langle \gamma_\tau F(w_\tau), w_\tau - u \rangle \\ &\leq \Theta(r_0) + \sum_{\tau=1}^t \frac{3\gamma_\tau^2}{2\alpha} [M^2 + (\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2] + \sum_{\tau=1}^t \langle \gamma_\tau \Delta_\tau, w_\tau - u \rangle \\ (5.11) \quad &= \Theta(r_0) + \sum_{\tau=1}^t \frac{3\gamma_\tau^2}{2\alpha} [M^2 + (\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2] + \sum_{\tau=1}^t \langle \gamma_\tau \Delta_\tau, w_\tau - y_{\tau-1} \rangle \\ &\quad + \sum_{\tau=1}^t \langle \gamma_\tau \Delta_\tau, y_{\tau-1} - u \rangle, \end{aligned}$$

where y_τ are given by (3.3). Since the sequences $\{y_\tau\}$, $\{\xi_\tau = \gamma_\tau \Delta_\tau\}$ satisfy the premise of Corollary 5.2, we have

$$\begin{aligned} (\forall u \in Z) : \sum_{\tau=1}^t \langle \gamma_\tau \Delta_\tau, y_{\tau-1} - u \rangle &\leq V(r_0, u) + \sum_{\tau=1}^t \frac{\gamma_\tau^2}{2\alpha} \|\Delta_\tau\|_*^2 \\ &\leq \Theta(r_0) + \sum_{\tau=1}^t \frac{\gamma_\tau^2}{2\alpha} \epsilon_{w_\tau}^2, \end{aligned}$$

and thus (5.11) implies that for any $u \in Z$

$$\begin{aligned} (5.12) \quad \sum_{\tau=1}^t \langle \gamma_\tau F(w_\tau), w_\tau - u \rangle &\leq 2\Theta(r_0) + \sum_{\tau=1}^t \langle \gamma_\tau \Delta_\tau, w_\tau - y_{\tau-1} \rangle \\ &\quad + \sum_{\tau=1}^t \frac{3\gamma_\tau^2}{2\alpha} \left[M^2 + (\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2 + \frac{\epsilon_{w_\tau}^2}{3} \right] \end{aligned}$$

To complete the proof of (3.5) in the general case, note that since F is monotone, (5.12) implies that for all $u \in Z$,

$$\sum_{\tau=1}^t \gamma_\tau \langle F(u), w_\tau - u \rangle \leq \Gamma(t),$$

where

$$\Gamma(t) = 2\Theta(r_0) + \sum_{\tau=1}^t \frac{3\gamma_\tau^2}{2\alpha} \left[M^2 + (\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2 + \frac{\epsilon_{w_\tau}^2}{3} \right] + \sum_{\tau=1}^t \langle \gamma_\tau \Delta_\tau, w_\tau - y_{\tau-1} \rangle$$

(cf. (3.6)), whence

$$\forall (u \in Z) : \langle F(u), \widehat{z}_t - u \rangle \leq \left[\sum_{\tau=1}^t \gamma_\tau \right]^{-1} \Gamma(t).$$

When taking the supremum over $u \in Z$, we arrive at (3.5).

In the case of a Nash v.i., setting $w_\tau = (w_{\tau,1}, \dots, w_{\tau,m})$ and $u = (u_1, \dots, u_m)$ and recalling the origin of F , due to the convexity of $\phi_i(z_i, z^i)$ in z_i , for all $u \in Z$ we get from (5.12):

$$\sum_{\tau=1}^t \gamma_\tau \sum_{i=1}^m [\phi_i(w_\tau) - \phi_i(u_i, (w_\tau)^i)] \leq \sum_{\tau=1}^t \gamma_\tau \sum_{i=1}^m \langle F^i(w_\tau), (w_\tau)_i - u_i \rangle \leq \Gamma(t).$$

Setting $\phi(z) = \sum_{i=1}^m \phi_i(z)$, we get

$$\sum_{\tau=1}^t \gamma_\tau \left[\phi(w_\tau) - \sum_{i=1}^m \phi_i(u_i, (w_\tau)^i) \right] \leq \Gamma(t).$$

Recalling that $\phi(\cdot)$ is convex and $\phi_i(u_i, \cdot)$ are concave, $i = 1, \dots, m$, the latter inequality implies that

$$\left[\sum_{\tau=1}^t \gamma_\tau \right] \left[\phi(\hat{z}_t) - \sum_{i=1}^m \phi_i(u_i, (\hat{z}_t)^i) \right] \leq \Gamma(t),$$

or, which is the same,

$$\sum_{i=1}^m \left[\phi_i(\hat{z}_t) - \sum_{i=1}^m \phi_i(u_i, (\hat{z}_t)^i) \right] \leq \left[\sum_{\tau=1}^t \gamma_\tau \right]^{-1} \Gamma(t).$$

This relation holds true for all $u = (u_1, \dots, u_m) \in Z$; taking maximum of both sides in u , we get

$$\text{Err}_N(\hat{z}_t) \leq \left[\sum_{\tau=1}^t \gamma_\tau \right]^{-1} \Gamma(t). \quad \blacksquare$$

5.2. Proof of Theorem 3.5. In what follows, we use the notation from Theorem 3.2. By this theorem, in the case of constant stepsizes $\gamma_\tau \equiv \gamma$ we have

$$(5.13) \quad \text{Err}_{\text{vi}}(\hat{z}_t) \leq [t\gamma]^{-1} \Gamma(t),$$

where

$$(5.14) \quad \begin{aligned} \Gamma(t) &= 2\Theta + \frac{3\gamma^2}{2\alpha} \sum_{\tau=1}^t \left[M^2 + (\epsilon_{r_{\tau-1}} + \epsilon_{w_\tau})^2 + \frac{\epsilon_{w_\tau}^2}{3} \right] + \gamma \sum_{\tau=1}^t \langle \Delta_\tau, w_\tau - y_{\tau-1} \rangle \\ &\leq 2\Theta + \frac{7\gamma^2}{2\alpha} \sum_{\tau=1}^t \left[M^2 + \epsilon_{r_{\tau-1}}^2 + \epsilon_{w_\tau}^2 \right] + \gamma \sum_{\tau=1}^t \langle \Delta_\tau, w_\tau - y_{\tau-1} \rangle. \end{aligned}$$

For a Nash v.i., Err_{vi} in this relation can be replaced with Err_N .

Note that by description of the algorithm $r_{\tau-1}$ is a deterministic function of $\zeta^{N(\tau-1)}$ and w_τ is a deterministic function of $\zeta^{M(\tau)}$ for certain increasing sequences of integers $\{M(\tau)\}$, $\{N(\tau)\}$ such that $N(\tau-1) < M(\tau) < N(\tau)$. Therefore $\epsilon_{r_{\tau-1}}$ is a deterministic function of $\zeta^{N(\tau-1)+1}$, and ϵ_{w_τ} and Δ_τ are deterministic functions of

$\zeta^{M(\tau)+1}$. Denoting by \mathbf{E}_i the expectation w.r.t. ζ_i , we conclude that under assumption I we have

$$(5.15) \quad \mathbf{E}_{N(\tau-1)+1} \left\{ \epsilon_{r_{\tau-1}}^2 \right\} \leq M^2, \quad \mathbf{E}_{M(\tau)+1} \left\{ \epsilon_{w_\tau}^2 \right\} \leq M^2, \quad \|\mathbf{E}_{M(\tau)+1} \{\Delta_\tau\}\|_* \leq \mu,$$

and under assumption II, in addition,

$$(5.16) \quad \begin{aligned} \mathbf{E}_{N(\tau-1)+1} \left\{ \exp\{\epsilon_{r_{\tau-1}}^2 M^{-2}\} \right\} &\leq \exp\{1\}, \\ \mathbf{E}_{M(\tau)+1} \left\{ \exp\{\epsilon_{w_\tau}^2 M^{-2}\} \right\} &\leq \exp\{1\}. \end{aligned}$$

Now, let

$$\Gamma_0(t) = \frac{7\gamma^2}{2\alpha} \sum_{\tau=1}^t \left[M^2 + \epsilon_{r_{\tau-1}}^2 + \epsilon_{w_\tau}^2 \right].$$

We conclude by (5.15) that

$$(5.17) \quad \mathbf{E} \{ \Gamma_0(t) \} \leq \frac{21\gamma^2 M^2 t}{2\alpha}.$$

Further, $y_{\tau-1}$ clearly is a deterministic function of $\zeta^{M(\tau-1)+1}$, whence $w_\tau - y_{\tau-1}$ is a deterministic function of $\zeta^{M(\tau)}$. Therefore

$$(5.18) \quad \begin{aligned} \mathbf{E}_{M(\tau)+1} \{ \langle \Delta_\tau, w_\tau - y_{\tau-1} \rangle \} &= \langle \mathbf{E}_{M(\tau)+1} \{ \Delta_\tau \}, w_\tau - y_{\tau-1} \rangle \\ &\leq \mu \|w_\tau - y_{\tau-1}\| \leq 2\mu\Omega, \end{aligned}$$

where the concluding inequality follows from the fact that Z is contained in the $\|\cdot\|$ -ball of radius $\Omega = \sqrt{2\Theta/\alpha}$ centered at z_c , see (2.5). From (5.18) it follows that

$$\mathbf{E} \left\{ \gamma \sum_{\tau=1}^t \langle \Delta_\tau, w_\tau - y_{\tau-1} \rangle \right\} \leq 2\mu\gamma t\Omega.$$

Combining the latter relation, (5.13), (5.14) and (5.17), we arrive at (3.9). (i) is proved.

To prove (ii), observe, first, that setting

$$J_t = \sum_{\tau=1}^t \left[M^{-2} \epsilon_{r_{\tau-1}}^2 + M^{-2} \epsilon_{w_\tau}^2 \right],$$

we get

$$(5.19) \quad \Gamma_0(t) = \frac{7\gamma^2 M^2}{2\alpha} [t + J_t].$$

At the same time, we can write

$$J_t = \sum_{j=1}^{2t} \xi_j,$$

where $\xi_j \geq 0$ is a deterministic function of $\zeta^{I(j)}$ for certain increasing sequence of integers $\{I(j)\}$. Moreover, when denoting by \mathbf{E}_j conditional expectation over $\zeta^{I(j)}, \zeta^{I(j)+1}, \dots, \zeta^{I(j)-1}$ being fixed, we have

$$\mathbf{E}_j \{ \exp\{\xi_j\} \} \leq \exp\{1\},$$

see (5.16). It follows that

$$\begin{aligned}
\mathbf{E} \left\{ \exp \left\{ \sum_{j=1}^{k+1} \xi_j \right\} \right\} &= \mathbf{E} \left\{ \mathbf{E}_{k+1} \left\{ \exp \left\{ \sum_{j=1}^k \xi_j \right\} \exp \{ \xi_{k+1} \} \right\} \right\} \\
(5.20) \quad &= \mathbf{E} \left\{ \exp \left\{ \sum_{j=1}^k \xi_j \right\} \mathbf{E}_{k+1} \{ \exp \{ \xi_{k+1} \} \} \right\} \leq \exp \{ 1 \} \mathbf{E} \left\{ \exp \left\{ \sum_{j=1}^k \xi_j \right\} \right\}.
\end{aligned}$$

Whence $\mathbf{E}[\exp\{J\}] \leq \exp\{2t\}$, and applying the Tchebychev inequality, we get

$$\forall \Lambda > 0 : \text{Prob} \{ J > 2t + \Lambda t \} \leq \exp\{-\Lambda t\}.$$

Along with (5.19) it implies that

$$(5.21) \quad \forall \Lambda \geq 0 : \text{Prob} \left\{ \Gamma_0(t) > \frac{21\gamma^2 M^2 t}{2\alpha} + \Lambda \frac{7\gamma^2 M^2 t}{2\alpha} \right\} \leq \exp\{-\Lambda t\}.$$

Let now $\xi_\tau = \langle \Delta_\tau, w_\tau - y_{\tau-1} \rangle$. Recall that $w_\tau - y_{\tau+1}$ is a deterministic function of $\zeta^{M(\tau)}$. Besides this, we have seen that $\|w_\tau - y_{\tau-1}\| \leq D \equiv 2\Omega$. Taking into account (5.15), (5.16), we get

$$(5.22) \quad \begin{aligned} (a) \quad & \mathbf{E}_{M(\tau)+1} \{ \xi_\tau \} \leq \rho \equiv \mu D, \\ (b) \quad & \mathbf{E}_{M(\tau)+1} \{ \exp \{ \xi_\tau^2 R^{-2} \} \} \leq \exp \{ 1 \}, \text{ with } R = MD. \end{aligned}$$

Observe that $\exp\{x\} \leq x + \exp\{9x^2/16\}$ for all x . Thus (5.22.b) implies for $0 \leq s \leq \frac{4}{3R}$

$$(5.23) \quad \mathbf{E}_{M(\tau)+1} \{ \exp \{ s \xi_\tau \} \} \leq s\rho + \exp \{ 9s^2 R^2 / 16 \} \leq \exp \{ s\rho + 9s^2 R^2 / 16 \}.$$

Further, we have $s\xi_\tau \leq \frac{3}{8}s^2 R^2 + \frac{2}{3}\xi_\tau^2 R^{-2}$, hence for all $s \geq 0$,

$$\mathbf{E}_{M(\tau)+1} \{ \exp \{ s \xi_\tau \} \} \leq \exp \{ 3s^2 R^2 / 8 \} \mathbf{E}_{M(\tau)+1} \left\{ \exp \left\{ \frac{2\xi_\tau^2}{3R^2} \right\} \right\} \leq \exp \left\{ \frac{3s^2 R^2}{8} + \frac{2}{3} \right\}.$$

When $s \geq \frac{4}{3R}$, the latter quantity is $\leq 3s^2 R^2 / 4$, which combines with (5.23) to imply that for $s \geq 0$,

$$(5.24) \quad \mathbf{E}_{M(\tau)+1} \{ \exp \{ s \xi_\tau \} \} \leq \exp \{ s\rho + 3s^2 R^2 / 4 \}.$$

Acting as in (5.20), we derive from (5.24) that

$$s \geq 0 \Rightarrow \mathbf{E} \left\{ \exp \left\{ s \sum_{\tau=1}^t \xi_\tau \right\} \right\} \leq \exp \{ st\rho + 3s^2 t R^2 / 4 \},$$

and by the Tchebychev inequality, for all $\Lambda > 0$,

$$\text{Prob} \left\{ \sum_{\tau=1}^t \xi_\tau > t\rho + \Lambda R \sqrt{t} \right\} \leq \inf_{s \geq 0} \exp \{ 3s^2 t R^2 / 4 - s \Lambda R \sqrt{t} \} = \exp \{ -\Lambda^2 / 3 \}.$$

Finally, we arrive at

$$(5.25) \quad \text{Prob} \left\{ \gamma \sum_{\tau=1}^t \langle \Delta_\tau, w_\tau - y_{\tau-1} \rangle > 2\gamma \left[\mu t + \Lambda M \sqrt{t} \right] \Omega \right\} \leq \exp \{ -\Lambda^2 / 3 \}.$$

for all $\Lambda > 0$. Combining (5.13), (5.14), (5.21) and (5.25), we get (3.10). \blacksquare

5.3. Proof of Lemma 4.1.

Proof of (i). We clearly have $Z^o = X^o \times Y^o$, and $\omega(\cdot)$ is indeed continuously differentiable on this set. Let $z = (x, y)$ and $z' = (x', y')$, $z, z' \in Z$. Then

$$\begin{aligned} \langle \omega'(z) - \omega'(z'), z - z' \rangle &= \frac{1}{\alpha_x \Omega_x^2} \langle \omega'_x(x) - \omega'_x(x'), x - x' \rangle + \frac{1}{\alpha_y \Omega_y^2} \langle \omega'_y(y), y - y' \rangle \\ &\geq \frac{1}{\Omega_x^2} \|x - x'\|_x^2 + \frac{1}{\Omega_y^2} \|y - y'\|_y^2 \geq \|[x' - x; y' - y]\|^2. \end{aligned}$$

Thus, $\omega(\cdot)$ is strongly convex on Z , modulus $\alpha = 1$, w.r.t. the norm $\|\cdot\|$. Further, the minimizer of $\omega(\cdot)$ on Z clearly is $z_c = (x_c, y_c)$, and

$$\Theta = \frac{1}{\alpha_x \Omega_x^2} \Theta_x + \frac{1}{\alpha_y \Omega_y^2} \Theta_y = 1,$$

so that $\Theta = 1$, whence $\Omega = \sqrt{2\Theta/\alpha} = \sqrt{2}$.

Proof of (ii). $\mathbf{1}^0$. Let $z = (x, y)$ and $z' = (x', y')$ with $z, z' \in Z$. Observe that $\|y - y'\|_y \leq 2\Omega_y$ and thus

$$(5.26) \quad \|y'\|_y \leq 2\Omega_y$$

due to $0 \in Y$.

On the other hand, we have from (4.9) $F(z') - F(z) = [\Delta_x; \Delta_y]$, where

$$\begin{aligned} \Delta_x &= \sum_{\ell=1}^m [\phi'_\ell(x') - \phi'_\ell(x)]^* [\mathbf{A}_\ell^T y' + b_\ell] + \sum_{\ell=1}^m [\phi'_\ell(x)]^* \mathbf{A}_\ell [y' - y], \\ \Delta_y &= - \sum_{\ell=1}^m \mathbf{A}_\ell^* [\phi_\ell(x) - \phi_\ell(x')] + \Phi'_*(y') - \Phi'_*(y). \end{aligned}$$

We have

$$\begin{aligned} \|\Delta_x\|_{x,*} &= \max_{h \in \mathcal{X}} \max_{\|h\|_x \leq 1} \langle h, \sum_{\ell=1}^m [(\phi'_\ell(x') - \phi'_\ell(x))^* [\mathbf{A}_\ell^T y' + b_\ell] + [\phi'_\ell(x)]^* \mathbf{A}_\ell [y' - y]] \rangle_{\mathcal{X}} \\ &\leq \sum_{\ell=1}^m \left[\max_{\substack{h \in \mathcal{X} \\ \|h\|_x \leq 1}} \langle h, [\phi'_\ell(x') - \phi'_\ell(x)]^* [\mathbf{A}_\ell^T y' + b_\ell] \rangle_{\mathcal{X}} + \max_{h \in \mathcal{X}, \|h\|_x \leq 1} \langle h, [\phi'_\ell(x)]^* \mathbf{A}_\ell [y' - y] \rangle_{\mathcal{X}} \right] \\ &= \sum_{\ell=1}^m \left[\max_{\substack{h \in \mathcal{X} \\ \|h\|_x \leq 1}} \langle [\phi'_\ell(x') - \phi'_\ell(x)] h, \mathbf{A}_\ell^T y' + b_\ell \rangle_{\mathcal{X}} + \max_{h \in \mathcal{X}, \|h\|_x \leq 1} \langle [\phi'_\ell(x)] h, \mathbf{A}_\ell [y' - y] \rangle_{\mathcal{X}} \right] \\ &\leq \sum_{\ell=1}^m \left[\max_{h \in \mathcal{X}, \|h\|_x \leq 1} \|[\phi'_\ell(x') - \phi'_\ell(x)] h\|_{(\ell)} \|\mathbf{A}_\ell y' + b_\ell\|_{(\ell,*)} \right. \\ &\quad \left. + \max_{h \in \mathcal{X}, \|h\|_x \leq 1} \|\phi'_\ell(x) h\|_{(\ell)} \|\mathbf{A}_\ell [y' - y]\|_{(\ell,*)} \right]. \end{aligned}$$

Then by (4.2),

$$\begin{aligned} \|\Delta_x\|_{x,*} &\leq \sum_{\ell=1}^m \left[[L_x \|x - x'\|_x + M_x] [\|\mathbf{A}_\ell y'\|_{(\ell,*)} + \|b_\ell\|_{(\ell,*)}] + [L_x \Omega_x + M_x] \|\mathbf{A}_\ell [y' - y]\|_{(\ell,*)} \right] \\ &= [L_x \|x - x'\|_x + M_x] \sum_{\ell=1}^m [\|\mathbf{A}_\ell y'\|_{(\ell,*)} + \|b_\ell\|_{(\ell,*)}] + [L_x \Omega_x + M_x] \sum_{\ell=1}^m \|\mathbf{A}_\ell [y' - y]\|_{(\ell,*)} \\ &\leq [L_x \|x - x'\|_x + M_x] [\mathcal{A} \|y'\|_y + \mathcal{B}] + [L_x \Omega_x + M_x] \mathcal{A} \|y - y'\|_y, \end{aligned}$$

by definition of \mathcal{A} and \mathcal{B} . Next, due to (5.26) we get by definition of $\|\cdot\|$

$$\begin{aligned}\|\Delta_x\|_{x,*} &\leq [L_x\|x-x'\|_x + M_x][2\mathcal{A}\Omega_y + \mathcal{B}] + [L_x\Omega_x + M_x]\mathcal{A}\|y-y'\|_y \\ &\leq [L_x\Omega_x\|z-z'\| + M_x][2\mathcal{A}\Omega_y + \mathcal{B}] + [L_x\Omega_x + M_x]\mathcal{A}\Omega_y\|z-z'\|,\end{aligned}$$

what implies

$$(a) : \|\Delta_x\|_{x,*} \leq [\Omega_x[2\mathcal{A}\Omega_y + \mathcal{B}]L_x + 2\mathcal{A}\Omega_y[L_x\Omega_x + M_x]]\|z-z'\| + [2\mathcal{A}\Omega_y + \mathcal{B}]M_x$$

Further,

$$\begin{aligned}\|\Delta_y\|_{y,*} &= \max_{\eta \in \mathcal{Y}, \|\eta\|_y \leq 1} \langle \eta, -\sum_{\ell=1}^m \mathbf{A}_\ell^*[\phi_\ell(x) - \phi_\ell(x')] + \Phi'_*(y') - \Phi'_*(y) \rangle_{\mathcal{Y}} \\ &\leq \max_{\eta \in \mathcal{Y}, \|\eta\|_y \leq 1} \sum_{\ell=1}^m \langle \eta, \mathbf{A}_\ell^*[\phi_\ell(x) - \phi_\ell(x')] \rangle_{\mathcal{Y}} + \|\Phi'_*(y') - \Phi'_*(y)\|_{y,*} \\ &= \max_{\eta \in \mathcal{Y}, \|\eta\|_y \leq 1} \sum_{\ell=1}^m \langle \mathbf{A}_\ell \eta, \phi_\ell(x) - \phi_\ell(x') \rangle_{\mathcal{E}_\ell} + \|\Phi'_*(y') - \Phi'_*(y)\|_{y,*} \\ &\leq \max_{\eta \in \mathcal{Y}, \|\eta\|_y \leq 1} \sum_{\ell=1}^m \|\mathbf{A}_\ell \eta\|_{(\ell,*)} \|\phi_\ell(x) - \phi_\ell(x')\|_{(\ell)} \|\Phi'_*(y') - \Phi'_*(y)\|_{y,*} \\ &\leq \max_{\eta \in \mathcal{Y}, \|\eta\|_y \leq 1} \sum_{\ell=1}^m \|\mathbf{A}_\ell \eta\|_{(\ell,*)} [L_x\Omega_x + M_x] \|x-x'\|_x [L_y\|y-y'\|_y + M_y],\end{aligned}$$

by (4.2.b) and (4.5). Now

$$\|\Delta_y\|_{y,*} \leq \mathcal{A}[L_x\Omega_x + M_x]\|x-x'\|_x + [L_y\|y-y'\|_y + M_y],$$

and we come to

$$(b) : \|\Delta_y\|_{y,*} \leq [\Omega_x\mathcal{A}[L_x\Omega_x + M_x] + \Omega_y L_y] \|z-z'\| + M_y.$$

From (a) and (b) it follows that

$$\begin{aligned}\|F(z) - F(z')\|_* &\leq \Omega_x\|\Delta_x\|_{x,*} + \Omega_y\|\Delta_y\|_{y,*} \\ &\leq [\Omega_x^2[2\mathcal{A}\Omega_y + \mathcal{B}]L_x + 3\mathcal{A}\Omega_x\Omega_y[L_x\Omega_x + M_x] + L_y\Omega_y^2] \|z-z'\| \\ &\quad + \Omega_x[2\mathcal{A}\Omega_y + \mathcal{B}]M_x + \Omega_y M_y.\end{aligned}$$

We have justified (4.12)

2⁰. Let us verify (4.13). The first relation in (4.13) is readily given by (4.3.a,c). Let us fix $z = (x, y) \in Z$ and i , and let

$$\begin{aligned}(5.27) \quad \Delta &= F(z) - \Xi(z, \zeta_i) \\ &= \underbrace{\left[\sum_{\ell=1}^m [\phi'_\ell(x) - \mathbf{G}_\ell(x, \zeta_i)]^* [\mathbf{A}_\ell y + b_\ell] \right]}_{\Delta_x} \underbrace{- \sum_{\ell=1}^m \mathbf{A}_\ell^* [\phi_\ell(x) - f_\ell(x, \zeta_i)]}_{\Delta_y}.\end{aligned}$$

As we have seen,

$$(5.28) \quad \sum_{\ell=1}^m \|\psi_\ell\|_{(\ell,*)} \leq 2\mathcal{A}\Omega_y + \mathcal{B}$$

Besides this, for $u_\ell \in \mathcal{E}_\ell$ we have

$$\begin{aligned}
\left\| \sum_{\ell=1}^m \mathbf{A}_\ell^* u_\ell \right\|_{y,*} &= \max_{\eta \in \mathcal{Y}, \|\eta\|_y \leq 1} \left\langle \sum_{\ell=1}^m \mathbf{A}_\ell^* u_\ell, \eta \right\rangle_{\mathcal{Y}} = \max_{\eta \in \mathcal{Y}, \|\eta\|_y \leq 1} \left\langle \sum_{\ell=1}^m u_\ell, \mathbf{A}_\ell \eta \right\rangle_{\mathcal{Y}} \\
&\leq \max_{\eta \in \mathcal{Y}, \|\eta\|_y \leq 1} \left[\sum_{1 \leq \ell \leq m} \|u_\ell\|_{(\ell)} \|\mathbf{A}_\ell \eta\|_{(\ell,*)} \right] \\
(5.29) \quad &\leq \max_{\eta \in \mathcal{Y}, \|\eta\|_y \leq 1} \left[\max_{1 \leq \ell \leq m} \|u_\ell\|_{(\ell)} \right] \sum_{1 \leq \ell \leq m} \|\mathbf{A}_\ell \eta\|_{(\ell,*)} = \mathcal{A} \max_{1 \leq \ell \leq m} \|u_\ell\|_{(\ell)}.
\end{aligned}$$

Hence, setting $u_\ell = \phi_\ell(x) - f_\ell(x, \zeta_i)$ we obtain

$$(5.30) \quad \|\Delta_y\|_{y,*} = \left\| \sum_{\ell=1}^m \mathbf{A}_\ell^* [\phi_\ell(x) - f_\ell(x, \zeta)] \right\|_{y,*} \leq \mathcal{A} \underbrace{\max_{1 \leq \ell \leq m} \|\phi_\ell(x) - f_\ell(x, \zeta)\|_{(\ell)}}_{\xi = \xi(\zeta_i)}.$$

Further,

$$\begin{aligned}
\|\Delta_x\|_{x,*} &= \max_{h \in \mathcal{X}, \|h\|_x \leq 1} \left\langle h, \sum_{\ell=1}^m [\phi'_\ell(x) - \mathbf{G}_\ell(x, \zeta_i)]^* \psi_\ell \right\rangle_{\mathcal{X}} \\
&= \max_{h \in \mathcal{X}, \|h\|_x \leq 1} \sum_{\ell=1}^m \langle [\phi'_\ell(x) - \mathbf{G}_\ell(x, \zeta_i)] h, \psi_\ell \rangle_{\mathcal{X}} \\
&\leq \max_{h \in \mathcal{X}, \|h\|_x \leq 1} \sum_{\ell=1}^m \|[\phi'_\ell(x) - \mathbf{G}_\ell(x, \zeta_i)] h\|_{(\ell)} \|\psi_\ell\|_{(\ell,*)} \\
&\leq \sum_{\ell=1}^m \underbrace{\max_{h \in \mathcal{X}, \|h\|_x \leq 1} \|[\phi'_\ell(x) - \mathbf{G}_\ell(x, \zeta_i)] h\|_{(\ell)}}_{\xi_\ell = \xi_\ell(\zeta_i)} \underbrace{\|\psi_\ell\|_{(\ell,*)}}_{\rho_\ell}
\end{aligned}$$

Invoking (5.28), we conclude that

$$(5.31) \quad \|\Delta_x\|_{x,*} \leq \sum_{\ell=1}^m \rho_\ell \xi_\ell,$$

where all $\rho_\ell \geq 0$, $\sum_{\ell} \rho_\ell \leq 2\mathcal{A}\Omega_y + \mathcal{B}$ and

$$\xi_\ell = \xi_\ell(\zeta_i) = \max_{h \in \mathcal{X}, \|h\|_x \leq 1} \|[\phi'_\ell(x) - \mathbf{G}_\ell(x, \zeta_i)] h\|_{(\ell)}$$

Denoting by $p^2(\eta)$ the second moment of a scalar random variable η , observe that $p(\cdot)$ is a norm on the space of square summable random variables representable as deterministic functions of ζ_i , and that

$$p(\xi) \leq \Omega_x M_x, \quad p(\xi_\ell) \leq M_x$$

by (4.3.b,d). Now by (5.30), (5.31),

$$[\mathbf{E} \{ \|\Delta\|_*^2 \}]^{\frac{1}{2}} = [\mathbf{E} \{ \Omega_x^2 \|\Delta_x\|_{x,*}^2 + \Omega_y^2 \|\Delta_y\|_{y,*}^2 \}]^{\frac{1}{2}}$$

$$\begin{aligned}
&\leq p(\Omega_x \|\Delta_x\|_{x,*} + \Omega_y \|\Delta_y\|_{y,*}) \leq p\left(\Omega_x \sum_{\ell=1}^m \rho_\ell \xi_\ell + \Omega_y \mathcal{A}\xi\right) \\
&\leq \Omega_x \sum_{\ell} \rho_\ell \max_{\ell} p(\xi_\ell) + \Omega_y \mathcal{A}p(\xi) \\
&\leq \Omega_x [2\mathcal{A}\Omega_y + \mathcal{B}]M_x + \Omega_y \mathcal{A}\Omega_x M_x,
\end{aligned}$$

and the latter quantity is $\leq M$, see (4.12). We have established the second relation in (4.13).

3⁰. It remains to prove that in the case of (4.14), relation (4.15) takes place. To this end, one can repeat word by word the reasoning from item 2⁰ with the function $p_e(\eta) = \inf\{t > 0 : \mathbf{E}\{\exp\{\eta^2/t^2\}\} \leq \exp\{1\}\}$ in the role of $p(\eta)$. Note that similarly to $p(\cdot)$, $p_e(\cdot)$ is a norm on the space of random variables η which are deterministic functions of ζ_i and are such that $p_e(\eta) < \infty$. \blacksquare

5.4. Proof of Lemma 4.3. Item (i) can be verified exactly as in the case of Lemma 4.1; the facts expressed in (i) depend solely on the construction from Section 4.2 preceding the latter Lemma, and are independent of what are the setups for X, \mathcal{X} and Y, \mathcal{Y} .

Let us verify item (ii). Note that we are in the situation

$$(5.32) \quad \begin{aligned} \|(x, y)\| &= \sqrt{\|x\|_1^2/(2 \ln(n)) + \|y\|_1^2/(4 \ln(p^{(1)}))}, \\ \|(\xi, \eta)\|_* &= \sqrt{2 \ln(n)\|\xi\|_\infty^2 + 4 \ln(p^{(1)})\|\eta\|_\infty^2}. \end{aligned}$$

For $z = (x, y), z' = (x', y') \in Z$ we have

$$F(z) - F(z') = \left[\Delta_x = [\text{Tr}((y - y')A_1); \dots; \text{Tr}((y - y')A_n)]; \Delta_y = - \sum_{j=1}^n (x_j - x'_j)A_j \right].$$

whence

$$\begin{aligned}
\|\Delta_x\|_\infty &\leq \|y - y'\|_1 \max_{1 \leq j \leq n} |A_j|_\infty \leq \sqrt{2 \ln(n)} A_\infty \|z - z'\|, \\
\|\Delta_y\|_\infty &\leq \|x - x'\|_\infty \max_{1 \leq j \leq n} |A_j|_\infty \leq 2\sqrt{\ln(p^{(1)})} A_\infty \|z - z'\|,
\end{aligned}$$

and

$$\|(\Delta_x, \Delta_y)\|_* \leq [2 \ln(n) + 4 \ln(p^{(1)})] \|z - z'\|,$$

as required in (4.29). Further, relation (4.30.a) is clear from the construction of Ξ_k . To prove (4.30.b), observe that when $(x, y) \in Z$, we have (see (4.27), (4.28))

$$\|\Xi^x(x, y, \eta)\|_\infty \leq |\bar{y}_i| \max_{1 \leq j \leq n} |A_j^i|_\infty \leq A_\infty,$$

and, since $F^x(x, y) = \mathbf{E}\{\Xi^x(x, y, \zeta)\}$,

$$(5.33) \quad \|\Xi^x(x, y, \eta) - F^x(x, y)\|_\infty \leq 2A_\infty.$$

Clearly,

$$(5.34) \quad \|\Xi^y(x, y, \eta) - F^y(x, y)\|_\infty = |A_j - \sum_{j=1}^n x_j A_j|_\infty \leq 2A_\infty.$$

Applying [4, Theorem 2.1(iii), Example 3.2, Lemma 1], we derive from (5.33) and (5.34) that for every $(x, y) \in Z$ and every $i = 1, 2, \dots$ it holds

$$\mathbf{E} \left\{ \exp \left\{ \left\| \Xi_k^x(x, y, \zeta_i) - F^x(x, y) \right\|_\infty^2 / N_{k,x}^2 \right\} \right\} \leq \exp\{1\},$$

$$N_{k,x} = 2A_\infty \left(2 \exp\{1/2\} \sqrt{\ln(n)} + 3 \right) k^{-1/2}$$

and

$$\mathbf{E} \left\{ \exp \left\{ \left\| \Xi_k^y(x, y, \zeta_i) - F^y(x, y) \right\|_\infty^2 / N_{k,y}^2 \right\} \right\} \leq \exp\{1\},$$

$$N_{k,y} = 2A_\infty \left(2 \exp\{1/2\} \sqrt{\ln(p^{(1)})} + 3 \right) k^{-1/2}.$$

Combining the latter bounds with (5.32) we conclude (4.30.b). ■