



**HAL**  
open science

# Large-scale nonconvex optimization: randomization, gap estimation, and numerical resolution

Joseph Frédéric Bonnans, Kang Liu, Nadia Oudjane, Laurent Pfeiffer, Cheng Wan

► **To cite this version:**

Joseph Frédéric Bonnans, Kang Liu, Nadia Oudjane, Laurent Pfeiffer, Cheng Wan. Large-scale nonconvex optimization: randomization, gap estimation, and numerical resolution. 2023. hal-03631702v2

**HAL Id: hal-03631702**

**<https://hal.science/hal-03631702v2>**

Preprint submitted on 3 Feb 2023 (v2), last revised 16 Jun 2023 (v3)

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

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

# LARGE-SCALE NONCONVEX OPTIMIZATION: RANDOMIZATION, GAP ESTIMATION, AND NUMERICAL RESOLUTION\*

J. FRÉDÉRIC BONNANS<sup>†</sup>, KANG LIU<sup>‡</sup>, NADIA OUDJANE<sup>§</sup>, LAURENT PFEIFFER<sup>†</sup>,  
AND CHENG WAN<sup>§</sup>

**Abstract.** We address a large-scale and nonconvex optimization problem, involving an aggregative term. This term can be interpreted as the sum of the contributions of  $N$  agents to some common good, with  $N$  large. We investigate a relaxation of this problem, obtained by randomization. The relaxation gap is proved to converge to zeros as  $N$  goes to infinity, independently of the dimension of the aggregate. We propose a stochastic method to construct an approximate minimizer of the original problem, given an approximate solution of the randomized problem. McDiarmid’s concentration inequality is used to quantify the probability of success of the method. We consider the Frank-Wolfe (FW) algorithm for the resolution of the randomized problem. Each iteration of the algorithm requires to solve a subproblem which can be decomposed into  $N$  independent optimization problems. A sublinear convergence rate is obtained for the FW algorithm. In order to handle the memory overflow problem possibly caused by the FW algorithm, we propose a stochastic Frank-Wolfe (SFW) algorithm, which ensures the convergence in both expectation and probability senses. Numerical experiments on a mixed-integer quadratic program illustrate the efficiency of the method.

**Key words.** Large-scale and nonconvex optimization, aggregative optimization, relaxation, decentralization, Frank-Wolfe algorithm, concentration inequalities, multi-agent optimization, privacy-preserving methods.

**AMS subject classifications.** 49M20, 49M27, 90C06, 90C26

## 1. Introduction.

*Problem formulation.* This article is devoted to the theoretical analysis and the numerical resolution of the following large-scale, aggregative, and nonconvex optimization problem:

$$(P) \quad \inf_{x \in \mathcal{X}} J(x) := f(G(x)), \quad \text{where: } \begin{cases} G(x) = \frac{1}{N} \sum_{i=1}^N g_i(x_i) \\ \mathcal{X} = \prod_{i=1}^N \mathcal{X}_i. \end{cases}$$

Here,  $N$  can be seen as the number of agents and is assumed to be large. The main feature of this problem is the aggregative form of the function  $G$ , which is defined as the average of the  $N$  mappings  $g_i$ , each of which defined on some set  $\mathcal{X}_i$  with image in a real Hilbert space  $\mathcal{E}$ . These mappings are referred to as the contribution mappings. We will call  $G(x)$  the aggregate. Let  $q$  denote the dimension of  $\mathcal{E}$  (possibly  $q = +\infty$ ). While very few structural assumptions are made on the sets  $\mathcal{X}_i$  and the mappings  $g_i$ , we will typically assume that  $f$  is convex, with a Lipschitz-continuous gradient. A central idea in this work is that the problem can be well approximated by a convex problem when  $N$  is large.

---

\*J.F. Bonnans was partially supported by the FiME Lab Research Initiative (Institut Europlace de Finance). This article benefited from the support of the FMJH Program PGMO and from the support to this program from EDF.

<sup>†</sup> Université Paris-Saclay, CNRS, CentraleSupélec, Inria, Laboratoire des signaux et systèmes, 91190, Gif-sur-Yvette, France ([frederic.bonnans@inria.fr](mailto:frederic.bonnans@inria.fr), [laurent.pfeiffer@inria.fr](mailto:laurent.pfeiffer@inria.fr)).

<sup>‡</sup> Institut Polytechnique de Paris, CNRS, Ecole Polytechnique, CMAP, 91128 Palaiseau, France ([kang.liu@polytechnique.edu](mailto:kang.liu@polytechnique.edu)).

<sup>§</sup> OSIRIS Department, EDF Lab, Paris-Saclay, and FiME, France ([nadia.oudjane@edf.fr](mailto:nadia.oudjane@edf.fr), [cheng.wan.05@polytechnique.org](mailto:cheng.wan.05@polytechnique.org)).

In various examples of interest, the function  $f$  has a separable structure as defined below. It turns out that taking into account the separability of  $f$ , when possible, allows us to refine our theoretical results (more precisely, to reduce some of the constants of interest, see Remark 2.7). From now on, we suppose that  $\mathcal{E}$  is the Cartesian product of  $M$  separable Hilbert spaces denoted  $\mathcal{E}_j$ , for  $j = 1, \dots, M$ . We assume that  $f$  is additively separable, that is to say, we assume that

$$f(y) = \sum_{j=1}^M f_j(y_j), \quad \forall (y_1, \dots, y_M) \in \prod_{i=1}^M \mathcal{E}_j,$$

where  $f_j: \mathcal{E}_j \rightarrow \mathbb{R}$ . Note that when  $f$  is not separable, one can take  $M = 1$  and  $\mathcal{E}_1 = \mathcal{E}$ . We assume that the contribution mappings are of the form

$$g_i(x_i) = (g_{ij}(x_i))_{j=1, \dots, M}, \quad \text{where } g_{ij}: \mathcal{X}_i \rightarrow \mathcal{E}_j.$$

Hence the criterion  $J$  of problem (P) writes

$$(1.1) \quad J(x) = f(G(x)) = \sum_{j=1}^M f_j \left( \frac{1}{N} \sum_{i=1}^N g_{ij}(x_i) \right).$$

*Motivating examples.* A particularly interesting instance of (P) is the social welfare optimization problem investigated in a closely related paper by Mengdi Wang [49]. The cost function is the following:

$$(1.2) \quad \inf_{x_i \in \mathcal{X}_i} f_0 \left( \frac{1}{N} \sum_{i=1}^N h_i(x_i) \right) + \frac{1}{N} \sum_{i=1}^N l_i(x_i).$$

Following her terminology, the function  $h_i$  is the contribution of agent  $i$  to some common goods,  $f_0$  is a social cost function of the common goods, and  $l_i$  describes the individual preference of agent  $i$ . There are various applications fitting into the framework of (1.2), see [49]. In particular, some power system management problems can be modeled as (1.2). Such a problem is investigated in [41]:  $x_i$  represents the production profile of the generator  $i$ ,  $l_i(x_i)$  is its individual production cost,  $f_0$  denotes the demand elasticity or, equivalently, a penalty function that depends on the difference between the average production and some inflexible demand  $D$  (e.g.  $f_0 := \|\cdot - D\|^2$ ) so as to penalize the deviation of the overall production from the inflexible demand.

Let us also mention the *resource allocation problems*, investigated in [4], for example. These problems are of the form (1.2), where  $f_0$  is the indicatrix function of a given point  $y \in \mathcal{E}$ , modelling the resource to be allocated over the agents. These problems find applications in energy management (see for example [22] and [27]). They do not fit to the current framework but can be reasonably well approximated, replacing the indicatrix by a penalty function.

We present and discuss other examples in Section 5, arising from supervised learning and optimal control.

*Related works and methods.* Let us return to the general problem (P). Classical Lagrangian relaxation (Chapter XII of [26]) methods can be relevant here because the dual problem is separable in the sense below, thanks to the aggregative form of  $G$ . To see this, let us reformulate (P) as:  $\inf_{(x,v) \in \mathcal{X} \times \mathcal{E}} f(v)$ , subject to the constraint that  $v = G(x)$ . Its dual problem is:

$$(1.3) \quad \sup_{\lambda \in \mathcal{E}} \left( -f^*(\lambda) + \Phi(\lambda) \right),$$

where  $f^*$  is the Fenchel conjugate function of  $f$ , and  $\Phi(\lambda)$  is defined by

$$(1.4) \quad \Phi(\lambda) := \inf_{x \in \mathcal{X}} \langle \lambda, G(x) \rangle = \frac{1}{N} \sum_{i=1}^N \inf_{x_i \in \mathcal{X}_i} \langle \lambda, g_i(x_i) \rangle.$$

One sees that  $\Phi(\lambda)$  can be evaluated by solving  $N$  independent sub-problems, one for each  $i$  in  $\{1, \dots, N\}$ . Solving these sub-problems can be much easier than addressing frontally the original problem with  $N$  coupled variables. This approach has been extensively employed in convex settings [41, 39]. However, the nonconvexity of the problem raises two major difficulties: the potentially large duality gap and the reconstruction of a primal solution from the dual optimal solution.

These two difficulties are addressed by Wang in [49]. She proposed a convex relaxation of the problem, based on a geometrical approach, that allows to obtain an estimate of the duality gap of order  $\mathcal{O}(q^2/N^2)$ . Her main tool was the Shapley-Folkman lemma [45], which allows to show that the image of  $G$  is close to a convex set. This idea was already present in the seminal work of Aubin and Ekeland in [2], dealing with a different setting involving a coupling constraint. We refer the reader to [29] for the most recent improvements dealing with this class of problems. We also refer to [49] for a more exhaustive of mathematical works dedicated to the estimation of the duality gap, where a kind of convexification occurs. After having solved the dual problem by a cutting plane method and then found an approximate solution to the relaxed primal problem via a projection problem, Wang's method recovers an approximate solution to the original nonconvex problem, by computing a Shapley-Folkman decomposition of the aggregate with a standard linear programming approach.

There exist another important class of methods for large-scale optimization problems which are the block coordinate descent algorithm and its variants [6, 20]. These methods may not be applicable without additional assumptions on the sets  $\mathcal{X}_i$  and the maps  $g_i$  (in the current framework, the sets  $\mathcal{X}_i$  could be discrete). Even if we make additional regularity assumptions, they may be inefficient, in particular because the cost function  $J$  is not convex in general.

*Contributions and organization of the paper.* We first introduce in Section 2 a convex relaxation of the original problem (P). The relaxed problem is obtained by randomization, that is to say, we replace the variables  $x_i$  by probability measures  $\mu_i$  on  $\mathcal{X}_i$ . The contribution mappings  $g_i(x_i)$  are replaced by  $\int_{\mathcal{X}_i} g_i(x_i) d\mu_i(x_i)$ ; these terms are linear with respect to  $\mu_i$ . The resulting randomized cost function, denoted  $\mathcal{J}$ , is convex, and so is the randomized problem. We give a first upper bound of the relaxation gap of order  $\mathcal{O}(1/N)$ . The randomized problem has a stochastic interpretation: it amounts to replace the variables  $x_i$  by independent random variables  $X_i$  of probability distribution  $\mu_i$ , and to replace  $g_i(x_i)$  by the expectation of  $g_i(X_i)$ . To derive a good candidate (for (P)), given an approximate solution to the randomized problem  $\mu = (\mu_1, \dots, \mu_N)$ , we propose to simulate random variables  $X_i$  with probability distribution  $\mu_i$ . We will call this technique the *selection method*. We give a sharp estimate of the probability of error for the selection method. More precisely, we estimate the probability that  $J(X_1, \dots, X_N) \geq \mathcal{J}(\mu) + (\frac{C}{N} + \epsilon)$ , given  $\epsilon > 0$ . The proof relies on McDiarmid's inequality, a concentration inequality [35].

From a numerical point of view, our main contribution is a method which is parallelizable, which benefits from the convexity of the randomized problem, but avoids the difficulty of the manipulation of probability measures (arising in the formulation of the randomized problem). This could be achieved by combining the Frank-Wolfe (FW) algorithm [17, 28], applied to the randomized problem, and the selection method

described previously. The resulting algorithm, called stochastic Frank-Wolfe (SFW) algorithm, is described and analyzed in Section 3. Each iteration of the algorithm requires to solve a subproblem of the form (1.4), which is decomposable into  $N$  subproblems. Resorting to the selection method, we avoid to manipulate explicitly probability measures on the sets  $\mathcal{X}_i$ , which may otherwise cause memory issues. The SFW method is able to find an  $\mathcal{O}(1/N)$ -solution to problem  $\mathbf{P}$ . In addition, we estimate the probability that the iterate  $x_k$  is  $(\frac{C}{k} + \epsilon)$ -optimal, for  $k \leq 2N$ , where  $k$  is the iteration counter. This result relies on concentration inequalities for martingales [15] which generalize McDiarmid's inequality.

Let us note that many articles in the literature are dedicated to stochastic variants of the Frank-Wolfe algorithm. These variants are concerned with the situation where the cost function is in the form of the expectation of a random cost and where its gradient is evaluated by sampling. See for example [16, 24, 25, 38, 51], see also [18, 32] and the references therein. Let us emphasize that the stochasticity of our algorithm has another origin, namely the selection method. In all these articles, convergence is established in expectation; to our knowledge, only the article [46] quantifies the probability of success of some stochastic method based on the Frank-Wolfe algorithm.

Our last theoretical contribution is a sharp estimate of the relaxation gap, of order  $\mathcal{O}(q \wedge N/N^2)$ , where  $q$  is the (potentially infinite) dimension of the aggregate space  $\mathcal{E}$ . It is proved in Section 4. It relies on a geometrical relaxation of problem  $(\mathbf{P})$ , shown to be equivalent to the relaxation by randomization. The relaxation gap is estimated with the help of a measure of nonconvexity for sets (introduced in [10]) and with the help of the Shapley-Folkman lemma [45]. We also give an estimate of the price of decentralization (as defined by Wang in [49]). We conclude the section with a detailed comparison of our approach and the one of [49].

Section 5 is dedicated to examples and discussions on numerical aspects. We provide in section 6 numerical results for a mixed-integer linear-quadratic program.

### 1.1. Notations.

*On sets.* For two sets  $\mathcal{A}$  and  $\mathcal{B}$  in a normed vector space  $\mathcal{X}$ , we denote by  $d(\mathcal{A}) := \sup_{x,y \in \mathcal{A}} \|x-y\|_{\mathcal{X}}$  the diameter of  $\mathcal{A}$ , by  $\mathcal{A}+\mathcal{B} = \{x+y \mid x \in \mathcal{A}, y \in \mathcal{B}\}$  the Minkowski sum of  $\mathcal{A}$  and  $\mathcal{B}$ , by  $\lambda\mathcal{A} = \{\lambda x \mid x \in \mathcal{A}\}$  the scalar multiplication of  $\mathcal{A}$  with  $\lambda \in \mathbb{R}$  and by  $\text{conv}(\mathcal{A})$  the convex hull of  $\mathcal{A}$ . Note that  $\text{conv}(\mathcal{A} + \mathcal{B}) = \text{conv}(\mathcal{A}) + \text{conv}(\mathcal{B})$ .

For all  $i \in \{1, \dots, N\}$ , we denote  $\mathcal{X}_{-i} = \left(\prod_{i'=1}^{i-1} \mathcal{X}_{i'}\right) \times \left(\prod_{i'=i+1}^N \mathcal{X}_{i'}\right)$ . Given  $x \in \mathcal{X}$ , we denote  $x_{-i} = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N) \in \mathcal{X}_{-i}$ . From time to time, we represent  $x$  by the pair  $(x_i, x_{-i})$ .

*On functions.* Let  $\mathcal{H}$  be a real Hilbert space. Let  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  and  $\|\cdot\|_{\mathcal{H}}$  denote the corresponding scalar product and norm. Let  $F: \mathcal{H} \rightarrow \mathbb{R} \cup \{+\infty\}$ . The domain of  $F$ , denoted by  $\text{dom}(F)$ , is defined by  $\text{dom}(F) = \{x \mid F(x) \neq +\infty\}$ . When  $F$  is differentiable, we denote its gradient by  $\nabla F$ . The gradient is defined as a function from  $\mathcal{H}$  to itself. We say that  $\nabla F$  is  $L$ -Lipschitz on a subset  $\mathcal{A}$  of  $\mathcal{H}$  if for any  $x, y \in \mathcal{A}$ , we have

$$(1.5) \quad \|\nabla F(x) - \nabla F(y)\|_{\mathcal{H}} \leq L\|x - y\|_{\mathcal{H}}.$$

The subgradient of  $F$  at some point  $x \in \text{dom}(F)$  is denoted by  $\partial F(x)$  and defined by

$$\partial F(x) = \{p \in \mathcal{H} \mid F(y) \geq F(x) + \langle p, y - x \rangle, \forall y \in \mathcal{H}\}.$$

The Fenchel conjugate of  $F$  is denoted by  $F^*: H \rightarrow \mathbb{R}$  and defined by  $F^*(p) = \sup_{x \in \mathcal{H}} \langle p, x \rangle - F(x)$ .

*On measures.* Given a set  $\Omega$ , we denote by  $\delta_x$  the Dirac distribution at some point  $x \in \Omega$ . We denote by  $\mathcal{P}_\delta(\Omega)$  the set of finitely supported probability distributions, defined by

$$\mathcal{P}_\delta(\Omega) := \left\{ \sum_{k=1}^K \lambda_k \delta_{x_k} \mid K \in \mathbb{N}, (\lambda_k)_{k=1}^K \in (\mathbb{R}_+)^K, (x_k)_{k=1}^K \in \Omega^K, \sum_{k=1}^K \lambda_k = 1 \right\}.$$

Let  $\mu = \sum_{k=1}^K \lambda_k \delta_{x_k} \in \mathcal{P}_\delta(\Omega)$ . Given a Hilbert space  $\mathcal{H}$  and a mapping  $F: \Omega \rightarrow \mathcal{H}$ , we denote

$$E_\mu[F] = \sum_{k=1}^K \lambda_k F(x_k), \quad \sigma_\mu^2[F] = \sum_{k=1}^K \lambda_k \|F(x_k) - E_\mu[F]\|_{\mathcal{H}}^2.$$

In other words,  $E_\mu[F]$  is the integral of  $F$  with respect to the measure  $\mu$  and  $\sigma_\mu^2[F]$  is the variance of the probability measure  $\sum_{j=1}^J \lambda_j \delta_{F(x_j)}$ , in the sense of [48, Remark 7.5]. Finally, the Bernoulli distribution with parameter  $\omega \in [0, 1]$  is denoted by  $\text{Bern}(\omega)$ .

*On numbers and real-valued random variables.* We denote by  $m \wedge n$  the minimum of the numbers  $m$  and  $n$  in  $\mathbb{R} \cup \{+\infty\}$ . Let  $X$  be a real-valued random variable. The expectation of  $X$  is denoted by  $\mathbb{E}[X]$ , the variance of  $X$  is denoted by  $\text{Var}(X)$  and the conditional expectation of  $X$  w.r.t. some  $\sigma$ -algebra  $\mathcal{F}$  is denoted by  $\mathbb{E}[X \mid \mathcal{F}]$ . Given  $\mu \in \mathcal{P}_\delta(\Omega)$  and a random variable  $X$  in  $\Omega$ , the notation  $X \sim \mu$  indicates that  $\mu$  is the probability distribution of  $X$ .

**2. Relaxation by randomization and gap estimation.** In this section we first make a structural assumption on the general problem of interest, problem (P). Next we introduce a relaxation of the problem, obtained by randomization. We give an upper bound of the randomization gap in Proposition 2.6. Finally we propose a method to recover an approximate solution to (P), given an approximate solution to the randomized problem. Its performance is investigated in Theorem 2.10.

**2.1. Assumptions and constants.** We recall that  $\mathcal{E}$  is the Cartesian product of  $M$  separable real Hilbert spaces  $\mathcal{E}_j$ . We denote by  $\langle \cdot, \cdot \rangle_{\mathcal{E}_j}$  the associated scalar products and by  $\|\cdot\|_{\mathcal{E}_j}$  the corresponding norms. Let us emphasize that we will not consider any other norm in the spaces  $\mathcal{E}_j$ . We equip  $\mathcal{E}$  with the scalar product  $\langle \cdot, \cdot \rangle$ , defined by  $\langle (y_1, \dots, y_M), (y'_1, \dots, y'_M) \rangle = \sum_{j=1}^M \langle y_j, y'_j \rangle_{\mathcal{E}_j}$  and we denote by  $\|\cdot\|$  the corresponding norm.

For any  $i = 1, \dots, N$  and for any  $j = 1, \dots, M$ , we denote

$$S_{ij} := \{g_{ij}(x_i) \mid x_i \in \mathcal{X}_i\} \quad \text{and} \quad S_j := \frac{1}{N} \sum_{i=1}^N S_{ij}.$$

The following regularity assumption will be in force all along the article.

ASSUMPTION A. For  $i = 1, 2, \dots, N$  and  $j = 1, 2, \dots, M$ :

1. The range set  $S_{ij}$  in  $\mathcal{E}_j$  has finite diameter  $d_{ij} := d(S_{ij})$ .
2. The function  $f_j$  is  $L_j$ -Lipschitz on  $\text{conv}(S_j)$ .
3. The function  $f_j$  is continuously differentiable on a neighborhood of  $\text{conv}(S_j)$ , and  $\nabla f_j$  is  $\tilde{L}_j$ -Lipschitz on  $\text{conv}(S_j)$ , in the sense of (1.5).

We next define two constants  $C_0 > 0$  and  $C_1 > 0$  by

$$C_0 = \sum_{j=1}^M \left( L_j \max_{1 \leq i \leq N} \{d_{ij}\} \right), \quad \text{and} \quad C_1 = \frac{1}{N} \sum_{j=1}^M \left( \tilde{L}_j \sum_{i=1}^N d_{ij}^2 \right).$$

*Remark 2.1.* We will regularly employ notations of the form  $O(h(N, q, k))$ , where  $h$  is an explicit function of  $N$ ,  $q$  (the dimension of  $\mathcal{E}$ ), and  $k$  (some iteration counter). We use it to express the fact that some variable is bounded by  $Ch(N, q, k)$ , where the constant  $C$  only depends on  $(\max_{1 \leq i \leq N} d_{ij})_{j=1, \dots, M}$  and the Lipschitz moduli  $(L_j)_{j=1, \dots, M}$  and  $(\tilde{L}_j)_{j=1, \dots, M}$ . With this convention in mind, we have

$$C_0 = O(1) \quad \text{and} \quad C_1 = O(1).$$

*Remark 2.2.* Our results can be applied to aggregative problems of the form

$$\inf_{x \in \mathcal{X}} \sum_{j=1}^M f_j \left( \sum_{i=1}^N \hat{g}_{ij}(x_i) \right),$$

i.e. of the same form as in (P), but without the coefficient  $\frac{1}{N}$ . Indeed, it suffices to define  $g_{ij} = N\hat{g}_{ij}$  to come down to the formulation (P) and to use the fact that  $d(g_{ij}(\mathcal{X}_i)) = Nd(\hat{g}_{ij}(\mathcal{X}_i))$ . The introduction of the coefficient  $\frac{1}{N}$  induces a natural scaling of the problem as  $N$  increases. It also enables to us to highlight the convexification of the problem as  $N$  becomes large, assuming that the coefficients  $d_{ij}$  are uniformly bounded.

We state in the following lemma a straightforward inequality, exhibiting the role of the constant  $C_0$ . Note that the role of the constant  $C_1$  will be revealed in Lemma 2.6.

**LEMMA 2.3.** *Let Assumption A be satisfied. For all  $i \in \{1, \dots, N\}$ , for all  $x_{-i} \in \mathcal{X}_{-i}$ ,  $x_i$  and  $x'_i$  in  $\mathcal{X}_i$ , it holds:*

$$|J(x'_i, x_{-i}) - J(x_i, x_{-i})| \leq \frac{C_0}{N}.$$

**2.2. The randomized problem.** The *randomized problem* is obtained by replacing each optimization variable  $x_i$  by a probability measure  $\mu_i \in \mathcal{P}_\delta(\mathcal{X}_i)$ . The contribution mappings  $g_i(x_i)$  are replaced by their integral with respect to  $\mu_i$ ,  $E_{\mu_i}[g_i]$ . Denoting  $\mathcal{P}_\delta = \prod_{i=1}^N \mathcal{P}_\delta(\mathcal{X}_i)$ , we obtain

$$(PR) \quad \inf_{\mu \in \mathcal{P}_\delta} \mathcal{J}(\mu) := f \left( \frac{1}{N} \sum_{i=1}^N E_{\mu_i}[g_i] \right) = \sum_{j=1}^M f_j \left( \frac{1}{N} \sum_{i=1}^N E_{\mu_i}[g_{ij}] \right).$$

The following equality justifies the denomination of the relaxed problem: given  $\mu \in \mathcal{P}_\delta$  and given  $N$  random variables  $X_i$  in  $\mathcal{X}_i$  such that  $X_i \sim \mu_i$ , we have

$$(2.1) \quad \mathcal{J}(\mu) = f \left( \frac{1}{N} \sum_{i=1}^N \mathbb{E}[g_i(X_i)] \right).$$

*Remark 2.4.* Working with probability measures with finite support, we do not need to equip the sets  $\mathcal{X}_i$  with a topology and to consider regularity assumptions on the mappings  $g_i$ . Note that the original problem and the randomized one do not necessarily have a solution under the standing assumptions of the article.

Let  $J^*$  and  $\mathcal{J}^*$  denote the values of the primal problem (P) and the randomized problem (PR) respectively. One is interested in comparing  $J^*$  and  $\mathcal{J}^*$ . The next lemma gives a direct result for one direction of this comparison.

LEMMA 2.5. *Let Assumption A hold true. Then  $-\infty < \mathcal{J}^* \leq J^*$ .*

*Proof.* By the definitions of  $E_{\mu_i}[g_{ij}]$  and  $S_j$ , we have that  $\frac{1}{N} \sum_{i=1}^N E_{\mu_i}[g_{ij}] \in \text{conv}(S_j)$ . Since  $f_j$  is Lipschitz-continuous over the bounded set  $\text{conv}(S_j)$ , we deduce that  $\mathcal{J}^* > -\infty$ . Let  $x \in \mathcal{X}$ . Define  $\mu = (\delta_{x_1}, \dots, \delta_{x_N}) \in \mathcal{P}_\delta$ . Then  $\mathcal{J}(\mu) = J(x)$ . As a consequence, inequality  $\mathcal{J}^* \leq J^*$  follows.  $\square$

The *randomization gap* is then defined as

$$\text{randomization gap} = J^* - \mathcal{J}^* \geq 0.$$

Next we prove a first upper bound of the randomization gap, of order  $O(\frac{1}{N})$ .

PROPOSITION 2.6. *Let Assumption A hold true. Let  $\mu \in \mathcal{P}_\delta$  and let  $(X_i)_{i=1, \dots, N}$  denote  $N$  independent random variables such that  $X_i \sim \mu_i$ . Then,*

$$(2.2) \quad \mathbb{E}[J(X)] - \mathcal{J}(\mu) \leq \frac{1}{2N^2} \sum_{j=1}^M \left( \tilde{L}_j \sum_{i=1}^N \sigma_{\mu_i}^2 [g_{ij}] \right) \leq \frac{C_1}{2N},$$

where  $X = (X_1, \dots, X_N)$ . As a consequence,  $J^* - \mathcal{J}^* \leq \frac{C_1}{2N}$ .

*Proof.* Let us define  $Y_j = \frac{1}{N} (\sum_{i=1}^N g_{ij}(X_i))$ , for  $j = 1, \dots, M$ . Let us set  $Y = (Y_j)_{j=1, \dots, M}$ . We have

$$\mathbb{E}[J(X)] = \mathbb{E}[f(Y)] \quad \text{and} \quad \mathcal{J}(\mu) = f(\mathbb{E}[Y]).$$

Since the variables  $X_i$  are independent, the random variables  $g_{ij}(X_i)$  are also independent (for fixed  $j$ ). It follows that

$$\mathbb{E}[\|Y_j - \mathbb{E}[Y_j]\|_{\mathcal{E}_j}^2] = \frac{1}{N^2} \sum_{i=1}^N \mathbb{E}[\|g_{ij}(X_i) - \mathbb{E}[g_{ij}(X_i)]\|_{\mathcal{E}_j}^2] = \frac{1}{N^2} \sum_{i=1}^N \sigma_{\mu_i}^2 [g_{ij}].$$

By Assumption A, we have

$$f(Y) \leq f(\mathbb{E}[Y]) + \langle \nabla f(\mathbb{E}[Y]), Y - \mathbb{E}[Y] \rangle_{\mathcal{E}_j} + \frac{1}{2} \sum_{j=1}^M \left( \tilde{L}_j \|Y_j - \mathbb{E}[Y_j]\|_{\mathcal{E}_j}^2 \right).$$

Taking the expectation of the above inequality and recalling the definition of  $C_1$ , we deduce (2.2).  $\square$

Remark 2.7. As we explained in the introduction, our analysis covers the case of a non-separable cost  $f$  (when  $M = 1$ ), however, when  $f$  is separable, it is useful to take this property into account. The aim of this remark is to justify this fact. Let us assume (in this remark only) that  $f$  is indeed separable, i.e.  $M > 1$ . Let us treat  $f$  as a non-separable function. It is easy to verify that the mapping  $\nabla f$  is Lipschitz continuous with modulus  $(\max_{j=1, \dots, M} \tilde{L}_j)$ ; this estimate is tight. If we do not take into account the additive structure of  $f$  in the proof of Proposition 2.6, we end up with the following estimate:

$$\mathbb{E}[J(X)] \leq \mathcal{J}(\mu) + \frac{1}{2N^2} \left( \max_{j=1, \dots, M} \tilde{L}_j \right) \sum_{i=1}^N \sum_{j=1}^M \sigma_{\mu_i}^2 [g_{ij}],$$



which is less precise than inequality (2.2). The same kind of comment could be made for the constants appearing afterwards in the convergence results of our numerical method.

We finish this subsection with an equivalent relaxed problem in the situation when the sets  $\mathcal{X}_i$  (resp. the contribution functions  $g_i$ ) are identical. We refer to this situation as the *symmetric case*.

LEMMA 2.8. *Suppose that there exists a set  $\mathcal{X}$  and a function  $g: \mathcal{X} \rightarrow \mathcal{E}$  such that  $\mathcal{X}_i = \mathcal{X}$  and  $g_i = g$ , for all  $i$ . Then,*

$$(2.3) \quad \mathcal{J}^* = \inf_{\nu \in \mathcal{P}_\delta(\mathcal{X})} f(E_\nu[g]).$$

*Proof.* Let  $\nu \in \mathcal{P}_\delta(\mathcal{X})$ . Take  $\mu = (\nu, \dots, \nu) \in \mathcal{P}_\delta$ . It follows that  $f(E_\nu[g]) = \mathcal{J}(\mu)$ . As a consequence,  $\inf_{\nu \in \mathcal{P}_\delta(\mathcal{X})} f(E_\nu[g]) \leq \inf_{\mu \in \mathcal{P}_\delta} \mathcal{J}(\mu)$ . On the other hand, let  $\bar{\mu} = (\bar{\mu}_1, \dots, \bar{\mu}_N) \in \mathcal{P}_\delta$ . Take  $\bar{\nu} = \sum_{i=1}^N \bar{\mu}_i / N \in \mathcal{P}_\delta(\mathcal{X})$ . Then, we deduce that  $\mathcal{J}(\bar{\mu}) = f(E_{\bar{\nu}}[g])$ . The conclusion follows.  $\square$

The relaxed problem in (2.3) has a natural interpretation as a mean field relaxation: instead of considering an optimization problem with  $N$  symmetric agents, we consider an arbitrarily large number of agents and optimize their distribution  $\nu$ .

**2.3. Selection method.** Suppose that a minimizer or an approximate minimizer  $\mu$  of the randomized problem (PR) has been obtained. We address in this subsection the issue of recovering an approximate minimizer of the original problem (P) from  $\mu$ .

A naive approach would consist in *averaging* the measures  $\mu_i$ , assuming that the sets  $\mathcal{X}_i$  are convex. In such a case, one can define the point  $x_i = E_{\mu_i}[\text{Id}]$ . Another approach, motivated by Proposition 2.6, consists in sampling  $\mu$ , that is, in simulating  $N$  independent random variables  $(X_1, \dots, X_N)$ , with distributions  $X_i \sim \mu_i$ . This can be done without additional structural assumption on the sets  $\mathcal{X}_i$ , moreover, Proposition 2.6 ensures that for any  $\varepsilon > 0$ ,

$$(2.4) \quad \mathbb{P}\left[J(X_1, \dots, X_N) < \mathcal{J}(\mu) + \frac{C_1}{2N} + \varepsilon\right] > 0.$$

Of course, one can realize several samplings of  $\mu$  to increase the probability of finding a good candidate for the original problem. We will refer to this approach as the *selection method*.

Example 2.9. Consider the following instance of the problem (P), where  $N$  is a large even number:

$$(2.5) \quad \begin{cases} \text{minimize} & \left\{ J(x_1, x_2, \dots, x_N) = -\frac{1}{N} \sum_{i=1}^N x_i^2 + \left( \frac{1}{N} \sum_{i=1}^N x_i \right)^2 \right\}; \\ \text{subject to} & x_i \in [-1, 1], \quad i = 1, \dots, N. \end{cases}$$

It is easy to see that  $x^*$  is a minimizer of (2.5) if and only if  $x^*$  has  $N/2$  coordinates equal to 1 and the others equal to  $-1$ . In this example, the original and the relaxed problem have the same value,  $J^* = \mathcal{J}^* = -1$ . The relaxed problem does not have a unique solution. One of them is  $\tilde{\mu}_i = \frac{1}{2}(\delta_{-1} + \delta_1)$ . Averaging  $\tilde{\mu}$  as suggested above yields  $\hat{x} = (0, \dots, 0)$  and  $J(\hat{x}) = 0$ . Thus in this example, the averaging method yields a poor candidate, whatever the value of  $N$ .

On the other hand, the selection method yields good candidates when  $N$  is large. Indeed, assume that  $\mathbb{P}[X_i = -1] = \mathbb{P}[X_i = 1] = 1/2$ . When  $N$  is large, by the law of

large numbers [47], nearly half of the random variables  $X_i$  are equal to 1 while the others are equal to  $-1$ , with probability close to 1. Then in such a case  $X$  is almost a minimizer of (2.5).

The next theorem provides a sharp estimate of the probability in (2.4) and confirms the interest of the selection method for large values of  $N$ . It relies on a concentration inequality, *McDiarmid's inequality* [35], and its variant [15] (cf. Corollary A.2) of “variance type”. It is quite intuitive that if the probability measures  $\mu_i$  have a small variance (in a sense to be specified), then the selection method will be more efficient. The interest of taking into account the variances of the probability distributions will be revealed in the analysis of the stochastic Frank-Wolfe algorithm in Subsection 3.3.

**THEOREM 2.10.** *Let Assumption A be satisfied. Let  $\mu \in \mathcal{P}_\delta$  and let  $X_1, \dots, X_N$  be  $N$  independent random variables such that  $X_i \sim \mu_i$ . Let  $X = (X_1, \dots, X_N)$ . Then, for all  $\epsilon > 0$ ,*

$$(2.6) \quad \mathbb{P} \left[ J(X) < \mathcal{J}(\mu) + \frac{C_1}{2N} + \epsilon \right] \geq 1 - \exp \left( -\frac{2N\epsilon^2}{C_0^2} \right).$$

Assume further that for all  $i = 1, \dots, N$ , there exists a constant  $v_i$  such that

$$(2.7) \quad \sigma_{\mu_i}^2 [J(\cdot, x_{-i})] \leq v_i^2,$$

for all  $x_{-i} \in X_{-i}$ . Then (2.6) can be strengthened as:

$$(2.8) \quad \mathbb{P} \left[ J(X) < \mathcal{J}(\mu) + \sum_{j=1}^M \sum_{i=1}^N \frac{\tilde{L}_j}{2N^2} \sigma_{\mu_i}^2 [g_{ij}] + \epsilon \right] \geq 1 - \exp \left( -\frac{N\epsilon^2}{2 \left( \sum_{i=1}^N N v_i^2 + \frac{C_0\epsilon}{3} \right)} \right).$$

*Proof.* Combining Lemma 2.3 and McDiarmid's inequality [35], we obtain

$$\mathbb{P} [J(X) < \mathbb{E}[J(X)] + \epsilon] \geq 1 - \exp \left( -\frac{2N\epsilon^2}{C_0^2} \right).$$

Combining this estimate with the second inequality of Proposition 2.6, we obtain (2.6).

Estimate (2.8) is proved similarly, combining McDiarmid's inequality of “variance type” proved in Corollary A.2 and the first inequality of Proposition 2.6.  $\square$

We provide in the next lemma an explicit candidate for (2.7).

**LEMMA 2.11.** *Inequality (2.7) is satisfied with  $v_i^2 = \frac{1}{N^2} \left( \sum_{j=1}^M L_j^2 \right) \sigma_{\mu_i}^2 (g_i)$ .*

*Proof.* We first state a general following property: given a probability measure  $\mu$  and two maps  $h_1$  and  $h_2$  suitably defined, we have the inequality

$$(2.9) \quad \sigma_\mu^2 [h_1 \circ h_2] \leq L^2 \sigma_\mu^2 [h_2],$$

assuming that  $h_1$  is  $L$ -Lipschitz continuous. Let us prove this property. For any  $x$ , we have

$$\begin{aligned} \|h_1 \circ h_2(x) - E_\mu[h_1 \circ h_2]\|^2 &= \|h_1 \circ h_2(x) - h_1(E_\mu[h_2])\|^2 \\ &\quad + 2\langle h_1 \circ h_2(x) - h_1(E_\mu[h_2]), h_1(E_\mu[h_2]) - E_\mu[h_1 \circ h_2] \rangle \\ &\quad + \|h_1(E_\mu[h_2]) - E_\mu[h_1 \circ h_2]\|^2. \end{aligned}$$

Taking the expectation, we obtain that

$$\sigma_\mu^2[h_1 \circ h_2] = E_\mu \left[ \left\| h_1 \circ h_2 - h_1(E_\mu[h_2]) \right\|^2 \right] - \left\| h_1(E_\mu[h_2]) - E_\mu[h_1 \circ h_2] \right\|^2.$$

Since  $h_1$  is  $L$ -Lipschitz continuous, we have  $E_\mu \left[ \left\| h_1 \circ h_2 - h_1(E_\mu[h_2]) \right\|^2 \right] \leq L^2 \sigma_\mu[h_2]^2$ . Inequality (2.9) follows immediately. Next, it is easy to verify that the function  $f$  is  $L$ -Lipschitz continuous, with  $L = \left( \sum_{j=1}^M L_j^2 \right)^{1/2}$ . Using (2.9), we conclude that

$$\sigma_{\mu_i}^2[J(\cdot, x_{-i})] \leq L^2 \sigma_{\mu_i}^2 \left[ \frac{1}{N} g_i(\cdot) + C \right] = \frac{L^2}{N^2} \sigma_{\mu_i}^2[g_i],$$

where  $C = \frac{1}{N} \sum_{i' \neq i} g_{i'}(x_{i'})$  is regarded as a constant. The estimate follows.  $\square$

### 3. Stochastic Frank-Wolfe algorithm.

**3.1. Assumptions.** We introduce two new assumptions, which will be in force until the end of the article.

ASSUMPTION B. *For all  $j = 1, \dots, M$ , the function  $f_j: \mathcal{E}_j \rightarrow \mathbb{R}$  is convex over  $\text{conv}(S_j)$ .*

Let  $\mu^1$  and  $\mu^2$  lie in  $\mathcal{P}_\delta$ . Take  $\omega \in [0, 1]$ . Let  $\mu = (\mu_1, \dots, \mu_N)$  be defined, for any  $i = 1, \dots, N$ , by  $\mu_i = (1 - \omega)\mu_i^1 + \omega\mu_i^2$ . Here, the addition and the multiplication by a scalar are understood as usual in the set of signed measures. In the sequel, we simply denote  $\mu = (1 - \omega)\mu^1 + \omega\mu^2$ . We have  $\mu \in \mathcal{P}_\delta$ ; moreover,  $E_{\mu_i}[g_i] = (1 - \omega)E_{\mu_i^1}[g_i] + \omega E_{\mu_i^2}[g_i]$ , for any  $i = 1, \dots, N$ . Then, Assumption B implies that  $\mathcal{J}(\mu) \leq (1 - \omega)\mathcal{J}(\mu^1) + \omega\mathcal{J}(\mu^2)$ . In words, the randomized problem (PR) is convex.

In this section, we address the numerical resolution of the randomized problem (and the original problem) under Assumption B. Let us mention that this convexity assumption is natural for the application problems described in the introduction. It allows the application of the Frank-Wolfe algorithm (also called conditional gradient algorithm) [17], for which convergence can be established. The Frank-Wolfe algorithm requires to solve at each iteration a subproblem. Here, the subproblems can be decomposed in  $N$  optimization problems, which can be solved in parallel. This property is particularly interesting, since we aim at solving instances of (P) with large values of  $N$ . We do not detail here the practical resolution of the subproblems, which can only be investigated case by case. Instead, we make the following assumption. Let us set  $\mathcal{A} := \{\nabla f(y) \mid y \in \text{conv}(G(\mathcal{X}))\} \subset \mathcal{E}$ .

ASSUMPTION C. *For all  $i = 1, \dots, N$ , for all  $\lambda \in \mathcal{A}$ , the problem*

$$(3.1) \quad \inf_{x_i \in \mathcal{X}_i} \langle \lambda, g_i(x_i) \rangle$$

*has at least a solution. For all  $i = 1, \dots, N$ , we fix a map  $\mathbb{S}_i: \mathcal{A} \mapsto \mathcal{X}_i$  such that for any  $\lambda \in \mathcal{A}$ ,  $\mathbb{S}_i(\lambda)$  is a solution to (3.1).*

The map  $\mathbb{S}_i$  can be understood as a best-response function corresponding to agent  $i$ . The involved cost function is a linear combination of the contribution mappings  $g_{ij}$ , with  $j = 1, \dots, M$ . In problem (3.1),  $\lambda$  can be interpreted as a price variable associated with  $g_i(x_i)$ .

*Remark 3.1.* It is easy to find assumptions which ensure the existence of the map  $\mathbb{S}_i$ . For example, one can assume that  $\mathcal{X}_i$  is a compact set in a topological vector space and that  $g_i$  is continuous. Let us emphasize that Assumption C is essentially

an assumption of numerical nature:  $\mathbb{S}_i$  should be understood as the output of an (efficient) numerical procedure for the resolution of (3.1). The algorithms described afterwards largely rely on evaluations of  $\mathbb{S}_i$ .

**3.2. Basic Frank-Wolfe algorithm.** We first describe a rather direct application of the Frank-Wolfe algorithm, which is referred to as the basic Frank-Wolfe algorithm. The starting point of our numerical approach is the following lemma, the proof of which is straightforward.

LEMMA 3.2. *Let  $\lambda \in \mathcal{A}$  and let  $\bar{\mu} = (\bar{\mu}_1, \dots, \bar{\mu}_N) \in \mathcal{P}_\delta$ . Then,  $\bar{\mu}$  is a solution to*

$$(3.2) \quad \inf_{\mu \in \mathcal{P}_\delta} \left\langle \lambda, \frac{1}{N} \sum_{i=1}^N E_{\mu_i}(g_i) \right\rangle.$$

*if and only if for all  $i = 1, \dots, N$ ,  $\bar{\mu}_i$  is supported in  $\operatorname{argmin}_{x_i \in \mathcal{X}_i} \langle \lambda, g_i(x_i) \rangle$ .*

The cost function in (3.2) should be regarded as a linearization of  $\mathcal{J}$ , as needed in the abstract formulation of the Frank-Wolfe algorithm in [17]. An immediate consequence of Lemma 3.2 is that  $(\delta_{\mathbb{S}_1(\lambda)}, \dots, \delta_{\mathbb{S}_N(\lambda)})$  is a solution to (3.2). The resolution of problem (3.2) is a key step in the numerical procedures developed afterwards; let us emphasize that the maps  $\mathbb{S}_i(y)$  can be evaluated independently from each other, i.e. the resolution of (3.2) can be parallelized.

---

**Algorithm 1** Frank-Wolfe Algorithm

---

Initialization:  $\mu^0 \in \mathcal{P}_\delta$ .  
**for**  $k = 0, 1, \dots, K$  **do**  
    **Step 1: Resolution of the subproblems.**  
    Set  $y^k = \frac{1}{N} \sum_{i=1}^N E_{\mu_i^k}[g_i]$  and set  $\lambda^k = \nabla f(y^k)$ .  
    **for**  $i = 1, \dots, N$  **do**  
        Compute  $\bar{x}_i^k = \mathbb{S}_i(\lambda^k)$ .  
    **end for**  
    Set  $\bar{\mu}^k = (\delta_{\bar{x}_1^k}, \dots, \delta_{\bar{x}_N^k})$ .  
    **Step 2: Update.**  
    Set  $\omega_k = 2/(k+2)$ .  
    Set  $\mu^{k+1} = (1 - \omega_k)\mu^k + \omega_k \bar{\mu}^k$ .  
**end for**

---

The convergence analysis performed afterwards relies on standard arguments (compare our proof with [28]). We introduce the primal gap  $\gamma_k$  and the primal-dual gap  $\beta_k$ , defined by

$$(3.3) \quad \gamma_k = \mathcal{J}(\mu^k) - \mathcal{J}^*, \quad \beta_k = \langle \nabla f(y^k), y^k - \bar{y}^k \rangle, \quad \text{where: } \bar{y}^k = \frac{1}{N} \sum_{i=1}^N g_i(\bar{x}_i^k).$$

Note that  $\beta_k$  can be evaluated numerically. The following lemma shows that  $\beta_k$  is an upper bound of the primal gap  $\gamma_k$ .

LEMMA 3.3. *For all  $k \in \mathbb{N}$ ,  $\gamma_k \leq \beta_k$ .*

*Proof.* Let  $k \in \mathbb{N}$ . Let  $\mu \in \mathcal{P}_\delta$  and let  $y = \frac{1}{N} \sum_{i=1}^N E_{\mu_i}[g_i]$ . By Lemma 3.2, we have  $\langle \nabla f(y^k), \bar{y}^k \rangle \leq \langle \nabla f(y^k), y \rangle$ . Thus, using the convexity of  $f$ , we obtain

$$(3.4) \quad \beta_k = \langle \nabla f(y^k), y^k - \bar{y}^k \rangle \geq \langle \nabla f(y^k), y^k - y \rangle \geq f(y^k) - f(y) = \mathcal{J}(\mu^k) - \mathcal{J}(\mu).$$

Since  $\mu$  is arbitrary, we deduce that  $\beta_k \geq \mathcal{J}(\mu^k) - \mathcal{J}^* = \gamma_k$ .  $\square$

We have the following convergence result.

**PROPOSITION 3.4.** *Let Assumptions A, B, and C hold. Then, in Algorithm 1, for any  $K \in \mathbb{N}^*$ ,*

$$\gamma_K \leq \frac{2C_1}{K}.$$

*Proof.* As we will see, the result is a consequence of Lemma A.3, with  $C = \frac{C_1}{2}$  and  $u_k = 0$ . By Assumption A,

$$f(y^{k+1}) \leq f(y^k) + \langle \nabla f(y^k), y^{k+1} - y^k \rangle + \sum_{j=1}^M \frac{\tilde{L}_j}{2} \|y_j^{k+1} - y_j^k\|^2.$$

We have  $y^{k+1} - y^k = \omega_k(\bar{y}^k - y^k)$ . Therefore, by definition of  $\beta_k$ ,

$$(3.5) \quad f(y^{k+1}) \leq f(y^k) - \omega_k \beta_k + \omega_k^2 \sum_{j=1}^M \frac{\tilde{L}_j}{2} \|\bar{y}_j^k - y_j^k\|^2.$$

By definition,  $\|\bar{y}_j^k - y_j^k\|^2 = \frac{1}{N^2} \left\| \sum_{i=1}^N E_{\mu_i^k} [g_{ij}(\bar{x}_i^k) - g_{ij}(\cdot)] \right\|^2$ , thus by Cauchy-Schwarz inequality,

$$\|\bar{y}_j^k - y_j^k\|^2 \leq \frac{1}{N} \sum_{i=1}^N \|E_{\mu_i^k} [g_{ij}(\bar{x}_i^k) - g_{ij}(\cdot)]\|^2 \leq \frac{1}{N} \sum_{i=1}^N d_{ij}^2.$$

Combining the above estimate with (3.5) and using the inequality  $\gamma_k \leq \beta_k$  proved in Lemma 3.3, we obtain that  $\gamma_{k+1} \leq (1 - \omega_k)\gamma_k + \frac{C_1}{2}\omega_k^2$ . Thus Lemma A.3 applies, which concludes the proof.  $\square$

In the following remark, we give an alternative value of  $\omega_k$  in Step 2 of Algorithm 1, while preserving the convergence rate from the previous proposition.

*Remark 3.5.* For any  $k \in \mathbb{N}$ , denote  $h_k(\omega) = -\omega\beta_k + \frac{C_k}{2}\omega^2$ , where the constant  $C_k$  is defined by  $C_k = \sum_{j=1}^M \tilde{L}_j \|\bar{y}_j^k - y_j^k\|^2$ . In view of inequality (3.5), the result of Proposition 3.4 remains true if the sequence  $(\omega_k)_{k \in \mathbb{N}}$  is chosen such that for any  $k \in \mathbb{N}$ ,  $h(\omega_k) \leq h(\bar{\omega}_k)$ . The result remains in particular true for

$$(3.6) \quad \omega_k = \operatorname{argmin}_{\omega \in [0,1]} h(\omega) = \min \left( \frac{\beta_k}{C_k}, 1 \right).$$

The above proposition shows the convergence of the Frank-Wolfe algorithm. Yet the algorithm only provides a relaxed solution. In order to get a solution to the original problem, one can use the selection method introduced in Subsection 2.3. A first direct application of Proposition 2.6 yields the following. Let  $(X_1, \dots, X_N)$  be  $N$  independent random variables such that  $X_i \sim \mu_i^k$ , for all  $i$ . Then,

$$\mathbb{E}[J(X)] \leq J^* + \frac{2C_1}{k} + \frac{C_1}{2N}.$$

Therefore, from a theoretical point of view, there is no guaranty of improvements when  $k \gg N$  since, then, the error term  $\frac{2C_1}{k}$  becomes negligible in comparison with  $\frac{C_1}{2N}$ . The

following lemma provides a convergence result (in probability) for the combination of the Frank-Wolfe algorithm and the selection method, for a number of iterations  $k \leq N$ .

LEMMA 3.6. *Let  $(\mu_k)_{k \in \mathbb{N}}$  be the output of Algorithm 1. Let  $k \leq N$ . Let  $\zeta \in (0, 1)$ . Let  $n \in \mathbb{N}^*$  and let  $(X_i^j)_{i=1, \dots, N}^{j=1, \dots, n}$  be  $Nn$  independent random variables such that  $X_i^j \sim \mu_i^k$ . Let  $X^j = (X_1^j, \dots, X_N^j)$ . Then,*

$$(3.7) \quad \mathbb{P} \left[ \min_{j=1, \dots, n} J(X^j) < \mathcal{J}^* + \frac{3C_1}{k} \right] \geq 1 - \zeta, \quad \text{if } n \geq \frac{2C_0^2}{C_1^2} \frac{k^2}{N} \ln \left( \frac{1}{\zeta} \right).$$

*Proof.* Since  $k \leq N$ , we have  $\frac{C_1}{2N} \leq \frac{C_1}{2k}$ . Therefore, by Theorem 2.10,

$$\mathbb{P} \left[ \min_{j=1, \dots, n} J(X^j) < \mathcal{J}^* + \frac{2C_1}{k} + \frac{C_1}{2k} + \epsilon \right] \geq 1 - \exp \left( - \frac{2N\epsilon^2 n}{C_0^2} \right),$$

for any  $\epsilon > 0$ . Take  $\epsilon = \frac{C_1}{2k}$ . If  $n$  satisfies (3.7), then  $\exp \left( - \frac{2N\epsilon^2 n}{C_0^2} \right) \leq \zeta$ .  $\square$

**3.3. Stochastic Frank-Wolfe algorithm.** At each iteration of Algorithm 1, a new point  $\bar{x}_i^k$  is added to the support of each distribution  $\mu_i^k$ . Therefore, if at iteration  $K$ , the points  $(\bar{x}_i^k)_{k=0, \dots, K-1}$  are distinct from each other (for each  $i$ ), then  $KN$  places are needed to store the iterate  $\mu^K$ , which can be prohibitive as  $K$  becomes large. We propose in this subsection a variant of Algorithm 1 which significantly mitigates the risk of memory overflow. We call it the *Stochastic Frank-Wolfe* (SFW) algorithm, it is given in Algorithm 2 below.

---

**Algorithm 2** Stochastic Frank-Wolfe Algorithm

---

Initialization:  $x^0 \in \mathcal{X}$   
**for**  $k = 0, 1, 2, \dots, K$  **do**  
  **Step 1: Resolution of the subproblems.**  
  Compute  $y^k = \frac{1}{N} \sum_{i=1}^N g_i(x_i^k)$  and  $\lambda^k = \nabla f(y^k)$ .  
  **for**  $i = 1, 2, \dots, N$  **do**  
    Compute  $\bar{x}_i^k = \mathbb{S}_i(\lambda^k)$ .  
  **end for**  
  **Step 2: Update.**  
  Choose  $n_k \in \mathbb{N}^*$ . Set  $\omega_k = 2/(k+2)$ .  
  **for**  $j = 1, 2, \dots, n_k$  **do**  
    **for**  $i = 1, 2, \dots, N$  **do**  
      Simulate  $P_i^{k,j} \sim \text{Bern}(\omega_k)$ , independently of all previously defined random variables.  
      Set  $\hat{x}_i^{k,j} = (1 - P_i^{k,j})x_i^k + P_i^{k,j}\bar{x}_i^k$ .  
    **end for**  
    Define  $\hat{x}^{k,j} = (\hat{x}_i^{k,j})_{i=1, \dots, N}$ .  
  **end for**  
  Find  $x^{k+1} \in \text{argmin} \{ J(x) \mid x \in \{\hat{x}^{k,j}, j = 1, 2, \dots, n_k\} \}$ .  
**end for**

---

Starting from an initialization  $x^0 \in \mathcal{X}$ , Algorithm 2 generates a sequence  $(x^k)_{k \in \mathbb{N}}$  in  $\mathcal{X}$ . Let us emphasize that there is no probability distribution involved in the practical implementation of Algorithm 2. However, for the analysis of the algorithm

and for its description, it is convenient to introduce  $\mu^k = (\delta_{x_1^k}, \dots, \delta_{x_N^k})$ . With this notation at hand, we first observe that  $y^k$ , as defined in Step 1 of Algorithm 2, satisfies  $y^k = \frac{1}{N} \sum_{i=1}^N E_{\mu_i^k}[g_i]$ . Thus the Steps 1 of Algorithms 1 and 2 play exactly the same role. Let us focus next on Step 2 of Algorithm 2 and let us define  $\bar{\mu}^k = (\delta_{\bar{x}_1^k}, \dots, \delta_{\bar{x}_N^k})$  and  $\hat{\mu}^k = (1 - \omega_k)\mu^k + \omega_k\bar{\mu}^k$ . In contrast with Algorithm 1, we do not directly use  $\hat{\mu}^k$  at the next iteration but instead employ our selection method so that  $\hat{\mu}^k$  is reduced to an  $N$ -uplet of Dirac measures. The application of the selection method is here simple since  $\hat{\mu}_i^k = (1 - \omega_k)\delta_{x_i^k} + \omega_k\delta_{\bar{x}_i^k}$ . Thus, to simulate a random variable with distribution  $\hat{\mu}_i^k$ , it suffices to simulate a random variable  $P$  with Bernoulli distribution  $\text{Bern}(\omega_k)$  and to consider  $(1 - P)x_i^k + P\bar{x}_i^k$ . Using this method, Step 2 consists in simulating  $n_k$  random variables  $(\hat{x}^{k,j})_{j=1, \dots, n_k}$  such that their probability distribution is equal to  $\hat{\mu}^k$  (to be rigorous, their probability distribution conditionally to  $x^k$ ). Finally, Step 2 selects a random variable  $\hat{x}^{k,j}$  which minimizes  $J$ .

It is important to keep in mind that all variables involved in the algorithm ( $x^k$ ,  $\bar{x}^k$ ,  $\hat{x}^{k,j}$ ) and all variables defined above ( $\mu^k$ ,  $\bar{\mu}^k$ ,  $\hat{\mu}^k$ ) are themselves random variables, since they depend on the Bernoulli random variables  $P_i^{k,j}$ . For the analysis of the algorithm, we need to consider the filtration generated by the Bernoulli random variables. We introduce the set of indices  $\mathcal{I}$  defined by

$$\mathcal{I} = \left\{ (k, j, i) \mid k \in \mathbb{N}, j \in \{1, \dots, n_k\}, i \in \{1, \dots, N\} \right\} \cup \{(0, 0, 0)\}.$$

We equip the set  $\mathcal{I}$  with the lexicographic order: given  $(k_1, j_1, i_1)$  and  $(k_2, j_2, i_2)$  in  $\mathcal{I}$ , we write  $(k_1, j_1, i_1) < (k_2, j_2, i_2)$  if and only if

$$[k_1 < k_2] \quad \text{or} \quad [k_1 = k_2 \text{ and } j_1 < j_2] \quad \text{or} \quad [(k_1, j_1) = (k_2, j_2) \text{ and } i_1 < i_2].$$

We further write  $(k_1, j_1, i_1) \leq (k_2, j_2, i_2)$  if and only if  $(k_1, j_1, i_1) < (k_2, j_2, i_2)$  or  $(k_1, j_1, i_1) = (k_2, j_2, i_2)$ . Note that this order coincides with the simulation order of the random variables  $P_i^{k,j}$  in the algorithm. The relation  $\leq$  defines a total order with minimal element  $(0, 0, 0)$ . For any  $(k, j, i) \neq (0, 0, 0)$ , we denote by  $(k, j, i) - 1$  the maximal element of the set  $\{(k', j', i') \in \mathcal{I} \mid (k', j', i') < (k, j, i)\}$ . Finally, we consider the filtration  $(\mathcal{G})_{(k,j,i) \in \mathcal{I}}$  defined by

$$\mathcal{G}_{(k,j,i)} = \begin{cases} \text{trivial } \sigma\text{-algebra,} & \text{if } (k, j, i) = (0, 0, 0), \\ \sigma(\mathcal{G}_{(k,j,i)-1}, P_i^{k,j}), & \text{otherwise,} \end{cases}$$

where  $\sigma(\mathcal{G}_{(k,j,i)-1}, P_i^{k,j})$  denotes the  $\sigma$ -algebra generated by  $\mathcal{G}_{(k,j,i)-1}$  and  $P_i^{k,j}$ . Note that  $\hat{x}_i^{k,j}$  is  $\mathcal{G}_{(k,j,i)}$ -adapted and that  $x^k$  and  $\bar{x}^k$  are  $\mathcal{G}_{(k,1,1)-1}$ -adapted.

**THEOREM 3.7.** *Let Assumptions A, B, and C hold true. Then, for all  $K = 1, \dots, 2N$ ,*

$$\mathbb{E}[\gamma_K] \leq \frac{4C_1}{K}, \quad \text{where } \gamma_K = J(x^K) - \mathcal{J}^*.$$

Moreover, for all  $\epsilon > 0$ ,

$$(3.8) \quad \mathbb{P}\left[\gamma_K < \frac{4C_1}{K} + \epsilon\right] \geq 1 - \exp\left(\frac{-\epsilon^2 N}{2(v_K + \epsilon m_K/3)}\right),$$

where  $v_K = \frac{2C_0^2}{K^2(K+1)^2} \left( \sum_{k=1}^{K-1} \frac{k(k+1)^2}{n_k} \right)$  and  $m_K = \frac{C_0}{K(K+1)} \left( \max_{k=1, \dots, K-1} \frac{(k+1)(k+2)}{n_k} \right)$ .

Finally, the following estimates quantify the variability of  $\gamma_K$ :

$$(3.9) \quad \text{Var}[\gamma_K] \leq \frac{16C_1^2}{K^2} + \frac{v_K}{N} \quad \text{and} \quad \mathbb{E}\left[\left(\max\left(\gamma_K - \frac{4C_1}{K}, 0\right)\right)^2\right] \leq \frac{v_K}{N}.$$

The proof is postponed to Section 3.4. Let us note that the constants  $m_K$  and  $v_K$  involved in the theorem depend on the sequence  $(n_k)_{k=0,1,\dots}$  but do not depend on  $N$ .

**COROLLARY 3.8.** *Let  $A > 0$ . Assume that  $n_k \geq \max\left(\frac{Ak^2}{N}, 1\right)$ , for any  $k$ . Then, for all  $K = 1, \dots, 2N$ ,*

$$\mathbb{P}\left[\gamma_K < \frac{4C_1 + C_0}{K}\right] \geq 1 - \exp\left(-\frac{A}{12}\right).$$

*Proof.* Using  $k + 1 \leq 2k$ , we obtain

$$\begin{aligned} v_K &\leq \frac{2C_0^2}{K^2(K+1)^2} \left( \sum_{k=1}^{K-1} \frac{Nk(k+1)^2}{Ak^2} \right) \leq \frac{8NC_0^2}{AK^2(K+1)^2} \left( \sum_{k=1}^{K-1} k \right) \\ &= \frac{4NC_0^2(K-1)K}{AK^2(K+1)^2} \leq \frac{4NC_0^2}{AK^2} \end{aligned}$$

and  $m_K \leq \frac{C_0}{K(K+1)} \left( \max_{k=1,\dots,K-1} \frac{N(k+1)(k+2)}{Ak^2} \right) \leq \frac{6NC_0}{AK^2}$ . Applying Theorem 3.7 with  $\epsilon = \frac{C_0}{K}$ , we obtain that  $\mathbb{P}\left[\gamma_K < \frac{4C_1 + C_0}{K}\right] \geq 1 - p$ , with

$$p \leq \exp\left(\frac{-(C_0/K)^2 N}{2\left(\frac{4NC_0^2}{AK^2} + \frac{6NC_0^2}{3AK^3}\right)}\right) \leq \exp\left(\frac{-A}{12}\right),$$

as was to be proved.  $\square$

*Remark 3.9.* A variant of Algorithm 2 consists in setting  $x^{k+1} = x^k$  if  $J(\hat{x}^{k,j}) \geq J(x^k)$  for all  $j = 1, \dots, n_k$ . Theorem 3.7 is still satisfied under this modification.

### 3.4. Proof of Theorem 3.7 and comments.

*Step 1: proof of the convergence in expectation.* We make use of the notations  $\mu^k$ ,  $\bar{\mu}^k$ , and  $\hat{\mu}^k$ , introduced right after Algorithm 2. We also introduce  $\beta_k = \langle \nabla f(y^k), y^k - \bar{y}^k \rangle$ , where  $\bar{y}^k = \frac{1}{N} \sum_{i=1}^N g_i(\bar{x}_i^k)$ . By construction, we have

$$J(x^{k+1}) = \min_{j=1,\dots,n_k} J(\hat{x}^{k,j}) \leq \frac{1}{n_k} \sum_{j=1}^{n_k} J(\hat{x}^{k,j}).$$

Recalling that  $\mathcal{J}(\mu^k) = J(x^k)$ , we deduce that  $\gamma_{k+1} \leq \gamma_k + a_k + b_k + c_k$ , where

$$\begin{aligned} a_k &= \frac{1}{n_k} \sum_{j=1}^{n_k} \left( J(\hat{x}^{k,j}) - \mathbb{E}[J(\hat{x}^{k,j}) | \mathcal{G}_{(k,1,1)-1}] \right), \\ b_k &= \frac{1}{n_k} \sum_{j=1}^{n_k} \left( \mathbb{E}[J(\hat{x}^{k,j}) | \mathcal{G}_{(k,1,1)-1}] - \mathcal{J}(\hat{\mu}^k) \right), \\ c_k &= \mathcal{J}(\hat{\mu}^k) - \mathcal{J}(\mu^k) = \mathcal{J}(\hat{\mu}^k) - J(x^k). \end{aligned}$$



The term  $a_k$  does not play a significant role at the moment since its expectation is null. The term  $b_k$  must be understood as a relaxation cost, induced by the use of the selection method. The term  $c_k$  is estimated exactly as in Proposition 3.4: as was seen in its proof, we have  $c_k \leq -\omega_k \beta_k + \omega_k^2 \frac{C_1}{2}$ . A direct adaptation of Proposition 2.6 shows that

$$b_k \leq \frac{1}{2N^2} \sum_{j=1}^M \sum_{i=1}^N \tilde{L}_j \sigma_{\tilde{\mu}_i}^2 [g_{ij}] \leq \frac{1}{2N^2} \sum_{j=1}^M \sum_{i=1}^N \tilde{L}_j \omega_k (1 - \omega_k) d_{ij}^2 = \omega_k (1 - \omega_k) \frac{C_1}{2N}.$$

Combining the above estimates, we obtain

$$(3.10) \quad \gamma_{k+1} \leq \gamma_k + a_k + \left( -\omega_k \beta_k + \omega_k^2 \frac{C_1}{2} \right) + \omega_k (1 - \omega_k) \frac{C_1}{2N}.$$

For the choice  $\omega_k = \bar{\omega}_k$ , we have  $(1 - \omega_k)/N = k/(N(k+2)) \leq \omega_k$ , since  $k \leq 2N$ . It follows that

$$\omega_k (1 - \omega_k) \frac{C_1}{2N} \leq \omega_k^2 \frac{C_1}{2}$$

and finally, since  $\gamma_k \leq \beta_k$ , we have  $\gamma_{k+1} \leq (1 - \omega_k) \gamma_k + \omega_k^2 C_1 + a_k$ . Next by Lemma A.3,

$$(3.11) \quad \gamma_K \leq \frac{4C_1}{K} + S_K, \quad \text{where: } S_K = \sum_{k=0}^{K-1} \frac{(k+1)(k+2)}{K(K+1)} a_k.$$

We have  $\mathbb{E}[a_k] = 0$ , thus  $\mathbb{E}[S_K] = 0$  and finally  $\mathbb{E}[\gamma_K] \leq \frac{4C_1}{K}$ .

*Step 2: proof of the probability and variance estimates.* We next need to find an estimate of  $\mathbb{P}[S_K \geq \epsilon]$ . For this purpose, we need to further decompose the term  $a_k$  as a sum of random variables. A first observation is the following equality:  $\mathbb{E}[J(\hat{x}^{k,j}) \mid \mathcal{G}_{(k,1,1)-1}] = \mathbb{E}[J(\hat{x}^{k,j}) \mid \mathcal{G}_{(k,j,1)-1}]$ , which easily follows from Lemma A.5. As a consequence,

$$J(\hat{x}^{k,j}) - \mathbb{E}[J(\hat{x}^{k,j}) \mid \mathcal{G}_{(k,1,1)-1}] = \sum_{i=1}^N U_{(k,j,i)},$$

where

$$U_{(k,j,i)} = \mathbb{E}[J(\hat{x}^{k,j}) \mid \mathcal{G}_{(k,j,i)}] - \mathbb{E}[J(\hat{x}^{k,j}) \mid \mathcal{G}_{(k,j,i)-1}].$$

We obtain the following decomposition of  $S_K$ :

$$S_K = \sum_{k=1}^{K-1} \sum_{j=1}^{n_k} \sum_{i=1}^N \frac{(k+1)(k+2)}{n_k K(K+1)} U_{(k,j,i)}.$$

Note that the index  $k$  starts at 1. Indeed,  $\omega_0 = 1$ , thus  $\hat{x}^{0,j} = \bar{x}^0$  and then  $a_0 = 0$ . Let us apply Proposition A.1 to  $S_K$ . We have  $\mathbb{E}[U_{(k,j,i)} \mid \mathcal{G}_{(k,j,i)-1}] = 0$ . Viewing the term  $J(\hat{x}^{k,j})$  as a function  $F$  of the random variables  $A := (P_{i'}^{k',j'})_{(k',j',i') < (k,j,i)}$ ,  $B := P_i^{k,j}$ , and  $C := (P_{i'}^{k',j'})_{(k,j,i) < (k',j',i')}$ , we can apply Lemma A.4 to  $U_{(k,j,i)}$ , with  $\delta = C_0/N$  (by Lemma 2.3). This yields

$$U_{(k,j,i)} \leq \frac{C_0}{N} \quad \text{and} \quad \mathbb{E}[U_{(k,j,i)}^2 \mid \mathcal{G}_{(k,j,i)-1}] \leq \frac{\omega_k (1 - \omega_k) C_0^2}{N^2}.$$

Therefore, Proposition A.1 applies to  $\mathbb{P}[S_K \geq \epsilon]$ , where the constants  $m$  and  $v$  are given by

$$m = \max_{k=1, \dots, K-1} \frac{(k+1)(k+2)C_0}{n_k K(K+1)} = \frac{m_K}{N},$$

$$v = \sum_{k=1}^{K-1} \sum_{j=1}^{n_k} \sum_{i=1}^N \left( \frac{(k+1)(k+2)}{n_k K(K+1)} \right)^2 \frac{2kC_0^2}{(k+2)^2 N^2} = \frac{v_K}{N}.$$

This proves estimate (3.8). Recalling that  $\gamma_K \leq \frac{4C_1}{K} + S_K$  a.s., we obtain

$$\text{Var}[\gamma_K] \leq \mathbb{E}[\gamma_K^2] \leq \mathbb{E}\left[\left(\frac{4C_1}{K} + S_K\right)^2\right] = \frac{16C_1^2}{K^2} + \mathbb{E}[S_K^2].$$

Next by Proposition A.1,  $\mathbb{E}[S_K^2] \leq v_K/N$ . The first inequality in (3.9) follows. The second inequality follows from the inequality:  $\max(\gamma_K - \frac{4C_1}{K}, 0)^2 \leq S_K^2$ .

*Remark 3.10.* Let us set  $h_k(\omega) = -\omega\beta_k + \omega^2\frac{C_1}{2} + \omega(1-\omega)\frac{C_1}{2N}$ . If for all  $k \in \mathbb{N}$ , we have  $h_k(\omega_k) \leq h_k(2/(k+2))$ , then the convergence in expectation of Theorem 3.7 still holds, i.e.  $\mathbb{E}[\gamma_K] \leq 4C_1/K$ , in view of inequality (3.10). In particular, one can take

$$(3.12) \quad \omega_k = \underset{\omega \in [0,1]}{\text{argmin}} h_k(\omega) = \max\left(\min\left(\frac{\beta_k - C_1/2N}{C_1(1-1/N)}, 1\right), 0\right).$$

**3.5. A speed-up of the SFW algorithm.** Step 1 of Algorithm 2 requires to solve  $N$  independent subproblems. It turns out that only a subset of those subproblems need to be solved for the implementation of Step 2. At iteration  $k$  consider the following set:

$$I_k = \bigcup_{j=1,2,\dots,n_k} \left\{ i \in \{1, \dots, N\} \mid P_i^{k,j} = 1 \right\}.$$

If  $i \notin I_k$ , then  $\hat{x}_i^{k,j} = x_i^k$ , in other words, for such an index  $i$ , it is not necessary to evaluate  $\mathbb{S}_i(\lambda^k)$ . A speed-up of the SFW algorithm can therefore be obtained by simulating the Bernoulli random variables before Step 1, next by evaluating  $\mathbb{S}_i(\lambda^k)$  only for the indices  $i$  in  $I_k$ , and finally by computing  $\hat{x}^{k,j}$  and  $x^{k+1}$  as before. The expectation of the number of subproblems to be solved at iteration  $k$  is given by

$$\begin{aligned} \mathbb{E}[|I_k|] &= \sum_{i=1}^N \mathbb{P}[i \in I_k] = N(1 - \mathbb{P}[1 \notin I_k]) = N\left(1 - \mathbb{P}[P_1^{k,j} = 0, \forall j = 1, \dots, n_k]\right) \\ &= N\left(1 - \left(\frac{k}{k+2}\right)^{n_k}\right). \end{aligned}$$

Note that this speed-up technique cannot be applied if  $\omega_k$  is chosen according to formula (3.12). Indeed, this formula requires to evaluate  $\beta_k$ , which implies that the  $N$  subproblems must all be solved.

**3.6. Stopping time strategy.** In Algorithm 2, the number of samplings  $n_k$  is chosen at the beginning of Step 2. We consider here a variant: we generate a sequence of random variables  $\hat{x}^{k,j}$  with probability distribution equal to  $\hat{\mu}_k$  (conditionally to

$\mathcal{G}_{(k,1,1)-1}$ ); the variables are constructed via Bernoulli variables independent from each other. We define  $n_k$  as the first index  $j$  such that

$$(3.13) \quad J(\hat{x}^{k,j}) \leq \mathcal{J}(\hat{\mu}^k) + \left(\frac{C_1}{2} + C_0\right)\omega_k^2,$$

or, equivalently,

$$(3.14) \quad f(\hat{y}^{k,j}) \leq f((1 - \omega_k)y^k + \omega_k\bar{y}^k) + \left(\frac{C_1}{2} + C_0\right)\omega_k^2,$$

where  $\bar{y}^k = \frac{1}{N} \sum_{i=1}^N g_i(\bar{x}_i^k)$  and  $\hat{y}^{k,j} = \frac{1}{N} \sum_{i=1}^N g_i(\hat{x}_i^{k,j})$ . The next iterate is defined by  $x^{k+1} = \hat{x}^{k,n_k}$ .

LEMMA 3.11. *Let  $(x^k)_{k \in \mathbb{N}}$  denote the sequence obtained with the stopping rule (3.13). Then*

$$J(x^{K+1}) - \mathcal{J}^* \leq \frac{4(C_1 + C_0)}{K}, \quad \forall K = 1, \dots, 2N, \quad a.s.$$

Moreover,

$$\mathbb{E}[n_k] \leq \left(1 - \exp\left(-\frac{4N}{(k+2)^3}\right)\right)^{-2}, \quad \forall k = 1, \dots, K.$$

*Proof.* Let  $\hat{x}$  be a random variable with probability distribution equal to  $\hat{\mu}^k$ , conditionally to  $\mathcal{G}_{(k,1,1)-1}$ . Then, for all  $\epsilon > 0$ , estimate (2.8) of Theorem 2.10 yields:

$$(3.15) \quad \mathbb{P}\left[J(\hat{x}) \geq \mathcal{J}(\hat{\mu}^k) + \frac{C_1}{2N}\omega_k(1 - \omega_k) + \epsilon \mid \mathcal{G}_{(k,1,1)-1}\right] \leq p_\epsilon$$

where  $p_\epsilon = \exp\left(\frac{-N\epsilon^2}{2(\omega_k(1-\omega_k)C_0^2 + \frac{C_0}{3}\epsilon)}\right)$ . For  $\epsilon = C_0\omega_k^2$ , we have

$$p_\epsilon = \exp\left(\frac{-NC_0^2\omega_k^4}{2(\omega_k C_0^2 - \frac{2}{3}\omega_k^2 C_0^2)}\right) \leq p := \exp\left(\frac{-N\omega_k^3}{2}\right) = \exp\left(\frac{-4N}{(k+2)^3}\right).$$

Recalling that  $\frac{C_1}{2N}\omega_k(1 - \omega_k) \leq \frac{C_1}{2}\omega_k^2$ , we deduce that

$$\mathbb{P}\left[J(\hat{x}) \geq \mathcal{J}(\hat{\mu}^k) + \left(\frac{C_1}{2} + C_0\right)\omega_k^2 \mid \mathcal{G}_{(k,1,1)-1}\right] \leq p.$$

Now, let us consider a sequence of independent random variables  $(\hat{x}^{k,j})_{j=1,\dots}$  (conditionally to  $\mathcal{G}_{(k,1,1)-1}$ ), with conditional probability distribution  $\hat{\mu}^k$ . By estimate (3.15),

$$\mathbb{P}[n_k = j] \leq \mathbb{P}\left[J(\hat{x}^{k,j'}) \geq \mathcal{J}(\hat{\mu}^k) + \left(\frac{C_1}{2} + C_0\right)\omega_k^2, \forall j' \mid \mathcal{G}_{(k,1,1)-1}\right] \leq p^{j-1}.$$

We finally deduce that  $\mathbb{E}[n_k] \leq \sum_{n=1}^{\infty} np^{n-1} = \frac{1}{(1-p)^2}$ , which proves the second part of the lemma. For the first part of the lemma, it suffices to observe that

$$J(x^{k+1}) \leq \mathcal{J}(\hat{\mu}^k) + \left(\frac{C_1}{2} + C_0\right)\omega_k^2 \leq J(x^k) - \beta_k\omega_k + (C_1 + C_0)\omega_k^2,$$

and to conclude with Lemma A.3.  $\square$

**3.7. Distributed algorithm.** In this subsection we present a privacy-preserving implementation of Algorithm 2. The Algorithm 3 is equivalent to Algorithm 2; the instructions are distributed over an **operator**,  $N$  **agents**, a **simulator**, and an **aggregator**, who communicate with each other. Roughly speaking, the operator sets up prices that are sent to the agents, which compute independently from each other their best-response. The aggregator computes in a confidential fashion the aggregate associated with a given value of  $(x_i)_{i=1,\dots,N}$ . The simulator implements the random variables  $P_i^{j,k}$  of the Stochastic Frank-Wolfe algorithm.

More precisely, at the beginning of iteration  $k$  of Algorithm 3, the operator sends a price  $\lambda_k$  to the agents, who calculate their best-response. The aggregator sends the corresponding aggregate  $\bar{y}_k$  to the operator, who can compute the primal-dual gap  $\beta^k$  and can fix the value of the stepsize  $\omega_k$ . Next the simulator realizes stochastic simulations, communicated to the agents. Only the aggregate associated with each simulation,  $\hat{y}^{k,j}$ , is communicated to the operator. The operator decides when to stop the simulation phase through the logical variable *test*. For example, *test* can be set to true as long as  $j < n_k$ , for predefined values of  $n_k$ . The variable *test* can also be designed so as to implement the stopping rule (3.14) of Subsection 3.6. Finally, the operator identifies the number  $j^*$  of the simulation that has yielded the best aggregate and communicates it to the agents.

The key point in this algorithm is that the operator never receives information that is specific to a given agent: it only collects aggregates (the variables  $\bar{y}^k$ ,  $\hat{y}^{k,j}$ , and  $y^k$ ). Similarly, the agents have only access to the prices  $\lambda_k$  and to  $j^*$ . We do not detail here algorithms used by the aggregator to compute the aggregate and refers the reader to [4], which investigates a similar approach for preserving privacy, with an operator that only has access to aggregates (note that the underlying mathematical method is different from ours). It is proposed in that reference to use a cryptographic protocol called *secure multiparty computation* for the non-intrusive computation of aggregates, taken from [43] and [1].

**Algorithm 3** Distributed SFW Algorithm

---

[**Agents**] Initialization:  $x^0 \in \mathcal{X}$ .  
 [**Aggregator**] Compute and send  $y^0 = \frac{1}{N} \sum_{i=1}^N g_i(x_i^0)$  to **Operator**.  
**for**  $k = 0, 1, 2, \dots, K$  **do**  
 [**Operator**] Compute and send  $\lambda^k = \nabla f(y^k)$  to the **Agents**.  
**for**  $i = 1, 2, \dots, N$  **do**  
 [**Agent**  $i$ ] Compute  $\bar{x}_i^k \in \mathbb{S}_i(\lambda^k)$ .  
**end for**  
 [**Aggregator**] Compute and send  $\bar{y}^k = \frac{1}{N} \sum_{i=1}^N g_i(\bar{x}_i^k)$  to **Operator**.  
 [**Operator**] Compute  $\beta^k = \langle \lambda^k, y^k - \bar{y}^k \rangle$ .  
 [**Operator**] Compute, send  $\omega_k$  with (3.12) or with  $\omega_k = \frac{2}{k+2}$  to **Simulator**.  
 [**Operator**] Set  $j = 0$  and send  $test = true$  to **Simulator**.  
**while**  $test$  **do**  
 [**Operator**] Increment  $j$ .  
**for**  $i = 1, 2, \dots, N$  **do**  
 [**Simulator**] Simulate and send  $P_i^{k,j} \sim \text{Bern}(\omega_k)$  to **Agent**  $i$ .  
 [**Agent**  $i$ ] Set  $\hat{x}_i^{k,j} = (1 - P_i^{k,j})x_i^k + P_i^{k,j}\bar{x}_i^k$ .  
**end for**  
 [**Aggregator**] Compute, send  $\hat{y}^{k,j} = \frac{1}{N} \sum_{i=1}^N g_i(\hat{x}_i^{k,j})$  to **Operator**.  
 [**Operator**] Update and send  $test$  to **Simulator**.  
**end while**  
 [**Operator**] Find  $j^* \in \underset{j'=1, \dots, j}{\text{argmin}} f(\hat{y}^{k,j'})$ . Set  $y^{k+1} = \hat{y}^{k,j^*}$ .  
 [**Operator**] Send  $j^*$  to the **Agents**.  
**for**  $i = 1, 2, \dots, N$  **do**  
 [**Agent**  $i$ ] Set  $x_i^{k+1} = \hat{x}_i^{k,j^*}$ .  
**end for**  
**end for**

---

**4. Refined gap estimates.**

**4.1. Nonconvexity measure and gap estimate.** We give in this subsection a refinement of the randomization gap obtained in Proposition 2.6. Our analysis relies on the concept of nonconvexity measure, introduced in [10].

DEFINITION 4.1. *Given a subset  $\mathcal{K}$  of  $\mathcal{E}$ , we call nonconvexity measure of  $\mathcal{K}$  the number  $\rho(\mathcal{K})$  defined by*

$$\rho(\mathcal{K}) = \left( \sup_{y \in \text{conv}(\mathcal{K})} \inf_{\substack{\mu \in \mathcal{P}_\delta, \\ E_\mu[\text{Id}] = y}} \sigma_\mu[\text{Id}]^2 \right)^{1/2},$$

where  $\text{Id}: \mathcal{E} \rightarrow \mathcal{E}$  denotes the identity mapping.

The "nonconvexity measure" terminology is motivated by the following: if  $\mathcal{K}$  is convex, then obviously  $\rho(\mathcal{K}) = 0$  and conversely, if  $\rho(\mathcal{K}) = 0$ , then  $\mathcal{K}$  is dense into  $\text{conv}(\mathcal{K})$ . We have the following two properties, easily verified. The map  $\rho$  is homogeneous in the following sense: given  $a \in \mathbb{R}$ , we have  $\rho(a\mathcal{K}) = |a|\rho(\mathcal{K})$ . Moreover  $\rho(\mathcal{K}) \leq d(\mathcal{K})$ , where  $d(\mathcal{K})$  is the diameter of  $\mathcal{K}$ . Another particularly interesting property for our aggregative problem is the sub-additivity of  $\rho(\cdot)^2$ : given two subsets  $\mathcal{K}_1$  and  $\mathcal{K}_2$ , we have  $\rho(\mathcal{K}_1 + \mathcal{K}_2)^2 \leq \rho(\mathcal{K}_1)^2 + \rho(\mathcal{K}_2)^2$ , see [10, Theorem 1]. We will

use an improvement of this inequality in the proof of Theorem 4.4, based on the Shapley-Folkman theorem.

The next lemma provides a general relaxation estimate based on the nonconvexity measure of the feasible set. Let us emphasize that the central idea behind this result is the same as the one in the proof of Proposition 2.6. The only difference is the point of view, which is here geometric while it was previously probabilistic.

LEMMA 4.2. *Let  $\mathcal{K}$  be a subset of  $\mathcal{E}$ . Let  $F$  be a differentiable real-valued function defined on some neighborhood of  $\text{conv}(\mathcal{K})$ . Assume that  $\nabla F$  is  $\tilde{L}$ -Lipschitz continuous over  $\text{conv}(\mathcal{K})$ . Then,*

$$\inf_{y \in \mathcal{K}} F(y) \leq \left( \inf_{y \in \text{conv}(\mathcal{K})} F(y) \right) + \frac{\tilde{L}}{2} \rho(\mathcal{K})^2.$$

*Proof.* Let  $y \in \text{conv}(\mathcal{K})$ . Let  $\mu \in \mathcal{P}_\delta(\mathcal{K})$  be such that  $E_\mu[\text{Id}] = y$ . Then, since  $\nabla F$  is  $\tilde{L}$ -Lipschitz continuous, we have

$$\inf_{y' \in \mathcal{K}} F(y') \leq E_\mu[F] \leq F(y) + \frac{\tilde{L}}{2} \sigma_\mu^2[\text{Id}].$$

Minimizing the right-hand side with respect to  $\mu$ , we obtain that

$$\inf_{y' \in \mathcal{K}} F(y') \leq F(y) + \frac{\tilde{L}}{2} \rho(\mathcal{K})^2.$$

Minimizing the result with respect to  $y$  yields the announced estimate.  $\square$

Some notations are needed for the application of Lemma 4.2 to (P). We set

$$\begin{aligned} \tilde{g}_{ij}(x_i) &= \sqrt{\tilde{L}_j} g_{ij}(x_i), & \tilde{g}_i(x_i) &= (\tilde{g}_{ij}(x_i))_{j=1, \dots, M} \\ \tilde{f}_j(y_j) &= f_j\left(\frac{y_j}{\sqrt{\tilde{L}_j}}\right), & \tilde{f}(y) &= \sum_{j=1}^M \tilde{f}_j(y_j). \end{aligned}$$

Obviously,  $J(x) = \tilde{f}\left(\frac{1}{N} \sum_{i=1}^N \tilde{g}_i(x_i)\right) = \sum_{j=1}^M \tilde{f}_j\left(\frac{1}{N} \sum_{i=1}^N \tilde{g}_{ij}(x_i)\right)$ . Finally we denote

$$\mathcal{Y}_i = \tilde{g}_i(\mathcal{X}_i) \quad \text{and} \quad \mathcal{Y} = \frac{1}{N} \sum_{i=1}^N \mathcal{Y}_i.$$

We give next two new formulations of problems (P) and (PR), revealing the geometric nature of the relaxation technique employed so far.

LEMMA 4.3. *We have*

$$(PG) \quad J^* = \inf_{y \in \mathcal{Y}} \tilde{f}(y),$$

$$(PGR) \quad \mathcal{J}^* = \inf_{y \in \text{conv}(\mathcal{Y})} \tilde{f}(y).$$

*Proof.* The first equality is straightforward. For the second one, it suffices to observe that  $\text{conv}(\mathcal{Y}) = \frac{1}{N} \sum_{i=1}^N \text{conv}(\mathcal{Y}_i)$  and that  $\text{conv}(\mathcal{Y}_i) = \{E_{\mu_i}[\tilde{g}_i] \mid \mu_i \in \mathcal{P}_\delta(\mathcal{X}_i)\}$ .  $\square$

We introduce the following constants:

$$(4.1) \quad D_i = \sum_{j=1}^M \tilde{L}_j d_{ij}^2, \quad D[k] = \max_{\substack{K \subseteq \{1, \dots, N\} \\ |K|=k}} \sum_{i \in K} D_i.$$

**THEOREM 4.4.** *Let Assumption A hold true. It holds:*

$$(4.2) \quad J^* - \mathcal{J}^* \leq \frac{1}{2N^2} \left( \max_{\substack{Q \subseteq \{1, \dots, N\} \\ |Q|=q \wedge N}} \sum_{i \in Q} \rho(\mathcal{Y}_i)^2 \right) \leq \frac{D[q \wedge N]}{2N^2}.$$

Note that  $D[N] = NC_1$ , thus the new gap estimate is the same as the one obtained in Proposition 2.6 when  $q \geq N$  and it is strictly better when  $q < N$ .

*Proof of Theorem 4.4.* We let the reader verify that  $\nabla \tilde{f}$  is 1-Lipschitz. Then Lemma 4.3 and the homogeneity of  $\rho$  yield

$$J^* - \mathcal{J}^* \leq \frac{1}{2} \rho(\mathcal{Y})^2 \leq \frac{1}{2N^2} \rho \left( \sum_{i=1}^N \mathcal{Y}_i \right)^2.$$

Applying [10, Theorem 2], we obtain that

$$\rho \left( \sum_{i=1}^N \mathcal{Y}_i \right)^2 \leq \max_{\substack{Q \subseteq \{1, \dots, N\} \\ |Q|=q \wedge N}} \sum_{i \in Q} \rho(\mathcal{Y}_i)^2,$$

which proves the first inequality. Observing that

$$\rho(\mathcal{Y}_i)^2 \leq d(\mathcal{Y}_i)^2 \leq \sum_{j=1}^M d(\tilde{g}_{ij}(\mathcal{X}_i))^2 = \sum_{j=1}^M \tilde{L}_j d(g_{ij}(\mathcal{X}_i))^2 = D_i,$$

we obtain the second inequality.  $\square$

**4.2. Duality and price of decentralization.** In this subsection we introduce a dual problem (we work again under Assumption B) and investigate its connection with the geometric relaxed problem (PGR). This allows us to obtain a last refinement of the randomization gap. For all  $i = 1, \dots, N$  and for all  $\lambda \in \mathcal{E}$ , we introduce

$$\Phi_i(\lambda) = \inf_{x_i \in \mathcal{X}_i} \langle \lambda, \tilde{g}_i(x_i) \rangle, \quad \mathcal{Y}_i(\lambda) = \operatorname{argmin}_{y_i \in \mathcal{Y}_i} \langle \lambda, y_i \rangle, \quad \mathcal{X}_i(\lambda) = \operatorname{argmin}_{x_i \in \mathcal{X}_i} \langle \lambda, \tilde{g}_i(x_i) \rangle.$$

We refer to the following problem as the dual problem:

$$(D) \quad \sup_{\lambda \in \mathcal{E}} \left( -\tilde{f}^*(\lambda) + \frac{1}{N} \sum_{i=1}^N \Phi_i(\lambda) \right).$$

Let  $\mathcal{D}^*$  denote the value of Problem (D).

**ASSUMPTION D.** *The function  $f: \mathcal{E} \rightarrow \mathbb{R}$  is lower semi-continuous and convex, and the set  $\operatorname{conv}(\mathcal{Y})$  is closed.*

*Remark 4.5.* Assume that  $\mathcal{E}$  is finite-dimensional. If the sets  $\mathcal{X}_i$  are compact and the maps  $\tilde{g}_i$  continuous, then the sets  $\mathcal{Y}_i = \tilde{g}_i(\mathcal{X}_i)$  are also compact. It is then easy to verify with Carathéodory's theorem that  $\operatorname{conv}(\mathcal{Y}_i)$  is also compact, thus closed, which finally implies Assumption D.

LEMMA 4.6. *The problem (PGR) has a solution.*

*Proof.* This is a direct application of [3, Theorem 11.9].  $\square$

The next lemma provides a duality result and a characterization of optimal solutions for problem (PGR).

LEMMA 4.7. *Let Assumptions A, B, C, and D hold true. Then,  $\mathcal{J}^* = \mathcal{D}^*$  and the dual problem (D) has at least one solution. Fix a solution  $\lambda$  to Problem (D). Let  $y \in \mathcal{E}$ . Then,  $y$  is a solution to (PGR) if and only if  $y \in \partial \tilde{f}^*(\lambda)$  and  $y \in \frac{1}{N} \sum_{i=1}^N \text{conv}(\mathcal{Y}_i(\lambda))$ .*

*Proof.* Let  $h$  denote the indicatrix function of  $\text{conv}(\mathcal{Y})$ . By Assumption A, the domain of  $\tilde{f}$  contains a neighborhood of  $\text{conv}(\mathcal{Y})$ . By Assumption D,  $h$  is lower semi-continuous. Therefore, the Fenchel-Rockafellar theorem [40] applies and yields

$$\mathcal{J}^* = \inf_{y \in \mathcal{E}} (f(y) + h(y)) = \sup_{\lambda \in \mathcal{E}} (-\tilde{f}^*(\lambda) - h^*(-\lambda)).$$

Moreover, the supremum in the right-hand side is a maximum. We have

$$-h^*(-\lambda) = \inf_{y \in \text{conv}(\mathcal{Y})} \langle \lambda, y \rangle = \inf_{y \in \mathcal{Y}} \langle \lambda, y \rangle = \frac{1}{N} \sum_{i=1}^N \Phi_i(\lambda).$$

As a consequence,  $\mathcal{J}^* = \mathcal{D}^*$  and problem (D) has at least one solution.

Now let us fix a solution  $\lambda$  to the dual problem (D). Let  $y \in \mathcal{E}$ . Then  $y$  is a solution if and only if (i)  $\tilde{f}(y) + \tilde{f}^*(\lambda) = \langle \lambda, y \rangle$  and (ii)  $h(y) + h^*(-\lambda) = -\langle \lambda, y \rangle$ . The condition (i) is equivalent to  $y \in \partial \tilde{f}(\lambda)$ . The condition (ii) is equivalent to

$$y \in \text{conv}(\mathcal{Y}) \text{ and } \langle \lambda, y \rangle = -h^*(-\lambda) = \inf_{y' \in \mathcal{Y}} \langle \lambda, y' \rangle.$$

Thus (ii)  $\iff y \in Y$ , where  $Y = \underset{y' \in \text{conv}(\mathcal{Y})}{\text{argmin}} \langle \lambda, y' \rangle$ . We further have

$$Y = \text{conv} \left( \underset{y' \in \mathcal{Y}}{\text{argmin}} \langle \lambda, y' \rangle \right) = \text{conv} \left( \frac{1}{N} \sum_{i=1}^N \mathcal{Y}_i(\lambda) \right) = \frac{1}{N} \sum_{i=1}^N \text{conv}(\mathcal{Y}_i(\lambda)),$$

which concludes the proof.  $\square$

*Remark 4.8.* If  $\tilde{f}$  is differentiable on  $\mathcal{E}$ , with a Lipschitz-continuous gradient, then  $\tilde{f}^*$  is strongly convex (see [3, Theorem 18.15]), which implies that (D) has a unique solution.

Let us fix a solution  $\lambda$  to the dual problem until the end of the subsection. Let us consider

$$J_{\text{dec}} = \inf_{x \in \mathcal{X}} J(x), \quad \text{subject to: } x_i \in \mathcal{X}_i(\lambda), \quad \forall i = 1, \dots, N.$$

In words, we restrict  $\mathcal{X}_i$  to the best-responses corresponding to the dual variable  $\lambda$ . Following the terminology of [49], we call price of decentralization the real number  $p = J_{\text{dec}} - \mathcal{J}^*$ .

PROPOSITION 4.9. *Let Assumptions A, B, C, and D hold true. It holds:*

$$p \leq J_{\text{dec}} - \mathcal{J}^* \leq \frac{1}{2N^2} \left( \max_{\substack{Q \subseteq \{1, \dots, N\} \\ |Q| = q \wedge N}} \sum_{i \in Q} \rho(\mathcal{Y}_i(\lambda))^2 \right).$$



*Proof.* The definition of  $J_{\text{dec}}$  and Lemma 4.7 respectively yield:

$$J_{\text{dec}} = \inf_{y \in \frac{1}{N} \sum_{i=1}^N \mathcal{Y}_i(\lambda)} \tilde{f}(y) \quad \text{and} \quad \mathcal{J}^* = \inf_{y \in \frac{1}{N} \sum_{i=1}^N \text{conv}(\mathcal{Y}_i(\lambda))} \tilde{f}(y).$$

The announced estimate follows then from Lemma 4.2 and [10, Theorem 2], as in the proof of Theorem 4.4.  $\square$

*Remark 4.10.* The randomization gap is bounded from above by  $J_{\text{dec}} - \mathcal{J}^*$ . Moreover, one can show that  $\rho(\mathcal{Y}_i(\lambda)) \leq \rho(\mathcal{Y}_i)$ . Thus Proposition 4.9 provides a last refinement of the gap estimate (4.2).

## 5. Comments on numerical aspects and examples.

**5.1. Literature comparison.** Let us compare our results and our method with the work of Wang [49]. Our gap estimate, as well as our estimate of the price of decentralization, are of order  $\mathcal{O}(\min(q, N)/N^2)$ , while the estimates obtained by applying [49, Theorem 3.5] are of order  $\mathcal{O}(q^2/N^2)$ . We emphasize that our first gap estimate, of order  $\mathcal{O}(1/N)$ , already improves [49] when  $q \gg \sqrt{N}$ . Note that the geometric relaxation employed in Section 4.1 is the same as the one used in [49].

Let us compare our algorithmic approaches. At a general level, one can observe that we have a primal approach, while Wang solves the dual problem to the relaxed problem. Our approach is restricted to the case where  $f$  is differentiable, while the dual approach allows to tackle the case of hard constraints (for example when  $f$  is the indicator function of some convex set). Both approaches leverage the decomposability of the problem into  $N$  problems and require that the subproblems can be easily solved. Let us emphasize however that we only need to be able to compute a single solution for those problems, while [49, Algorithm 2] requires to compute the full set of  $\xi$ -optimal solutions, which may be much more difficult. Our algorithm does not require to perform Shapley-Folkman decompositions, contrary to [49]. This is a major advantage when the dimension of the aggregate  $q$  is very large. Also, we do not need to evaluate  $f^*$ . As a counterpart, we are only able to find  $\mathcal{O}(1/N)$ -optimal solutions, while the algorithm of [49] can find  $\mathcal{O}(q^2/N^2)$ -optimal solutions. The design of a method for the computation of  $\mathcal{O}(q \wedge N/N^2)$ -solutions will be the topic of future research.

**5.2. Discussion.** The stochastic Frank-Wolfe algorithm investigated in the previous sections was motivated by the difficulty of manipulating probability measures, from a numerical point of view. However, when the sets  $\mathcal{X}_i$  are finite, with relatively low cardinality, it is possible to store probability measures with possibly full support and some other numerical methods can be used to solve the randomized problem. Let us assume (in this subsection only) that the sets  $\mathcal{X}_i$  are of cardinality  $n_i \in \mathbb{N}$  and that  $\mathcal{X}_i = \{\mathbf{x}_i^1, \dots, \mathbf{x}_i^{n_i}\}$ . Then the randomized problem reads:

$$(5.1) \quad \min_{\nu=(\nu_1, \dots, \nu_N)} f\left(\frac{1}{N} \sum_{i=1}^N \sum_{\ell=1}^{n_i} \nu_i^\ell g_i(\mathbf{x}_i^\ell)\right), \quad \text{subject to: } \nu_i \in \Delta(n_i),$$

where  $\Delta(n_i)$  denotes the  $(n_i - 1)$ -simplex, i.e.

$$\Delta(n_i) = \left\{ \nu \in \mathbb{R}^{n_i} \mid \sum_{\ell=1}^{n_i} \nu^\ell = 1 \text{ and } \nu^\ell \geq 0, \forall \ell = 1, \dots, n_i \right\}.$$

The problem is a convex program on a Cartesian product of  $N$  simplices. Let us first note that in this framework, Assumption C is trivially verified, since problem (3.1) is

just a minimization problem over  $\mathcal{X}_i$  which can be solved by enumeration. Moreover any variant of the Frank-Wolfe algorithm can be implemented, in order to solve the randomized problem in a faster way. We refer the reader to [28, 31]. Some other methods could also be implemented. The problem could be solved with the projected gradient descent algorithm, but the projection on the simplices is expensive (see [13]). Instead, the problem can be naturally addressed with the mirror descent algorithm [5] (see in particular the entropic descent algorithm in Section 5), and with accelerated versions of the entropic descent algorithm [30].

Let us observe that if we require  $\nu$  to have integer entries in the problem (5.1), then we are back to the original problem. Indeed, the elements of the simplex with integer entries are its vertices, that is, the vectors of the form  $(0, \dots, 0, 1, 0, \dots, 0)$ . Therefore the original problem can be viewed as a mixed-integer convex program (MICEP) and can be addressed numerically with combinatorial techniques, see [8, 12] and the references therein.

**5.3. Aggregative optimal control.** We describe here a large-scale optimal control problem of the form of problem (P), with an infinite-dimensional aggregate space. We verify Assumptions A, B, and C and we discuss the applicability of the Stochastic Frank-Wolfe algorithm.

Let us first fix the data of the problem. For any  $i = 1, \dots, N$ , we consider: an initial condition  $z_i^0 \in \mathbb{R}^n$ , a control set  $U_i \subseteq \mathbb{R}^m$ , a dynamics  $F_i: (z_i, u_i) \in \mathbb{R}^n \times U_i \mapsto F_i(z_i, u_i) \in \mathbb{R}^n$ , and a contribution function  $\phi_i: \mathbb{R}^n \times U_i \rightarrow \mathbb{R}^k$ . We also consider a social cost  $\ell: \mathbb{R}^k \rightarrow \mathbb{R}$ . We make the following assumptions:

1. *Regularity and boundedness.* For any  $i = 1, \dots, N$ ,
  - $U_i$  is non-empty and compact
  - $F_i$  is continuous, Lipschitz continuous with respect to  $z_i$ , uniformly with respect to  $u_i$ ; moreover, there exists a constant  $K_i$  such that  $\|F_i(z_i, u_i)\| \leq K_i(1 + \|z_i\|)$ , for any  $(z_i, u_i) \in \mathbb{R}^n \times U_i$
  - $\phi_i$  is continuous; moreover, there exists a function  $R_i: \mathbb{R}_+ \rightarrow \mathbb{R}_+$  such that  $\|\phi_i(z_i, u_i)\| \leq R_i(\|z_i\| + \|u_i\|)$ , for any  $(z_i, u_i) \in \mathbb{R}^n \times U_i$ .
2. *Regularity of the social cost.* The function  $\ell$  is continuously differentiable, moreover,  $\ell$  and  $\nabla \ell$  are Lipschitz continuous with moduli  $L_\ell$  and  $L_{\nabla \ell}$ , respectively.
3. *Convexity assumption.* For any  $i = 1, \dots, N$ , for any  $y \in \mathbb{R}^k$ , for any  $z_i \in \mathbb{R}^n$ , we define  $\mathcal{Z}_i(y, z_i)$  the set of all elements  $(\bar{z}_1, \bar{z}_2)$  in  $\mathbb{R}^{n+1}$ , where there exists  $u_i \in U_i$ , such that  $\bar{z}_1 = F_i(z_i, u_i)$  and  $\bar{z}_2 \geq \langle \nabla \ell(y), \phi_i(z_i, u_i) \rangle$ . The set  $\mathcal{Z}_i(y, z_i)$  is convex.

Let us mention a particular case in which the above convexity assumption is true: for any  $i = 1, \dots, N$ , for any  $y \in \mathbb{R}^k$ , for any  $z_i \in \mathbb{R}^n$ ,

- For any  $z_i$ , the map  $u_i \mapsto F_i(z_i, u_i)$  is affine.
- The set  $U_i$  is convex and the function  $u_i \in U_i \mapsto \langle \nabla \ell(y), \phi_i(z_i, u_i) \rangle$  is convex.

For any  $i = 1, \dots, N$ , consider the set  $\mathcal{X}_i$  of pairs  $(z_i, u_i) \in W^{1,\infty}(0, T; \mathbb{R}^n) \times L^\infty(0, T; \mathbb{R}^m)$  satisfying

$$\dot{z}_i(t) = F_i(z_i(t), u_i(t)), \quad z_i(0) = z_i^0, \quad u_i(t) \in U_i, \quad \text{for a.e. } t \in (0, T).$$

A direct application of Gronwall's lemma shows that for any  $(z_i, u_i) \in \mathcal{X}_i$ , we have  $\|z_i\|_{L^\infty(0, T; \mathbb{R}^n)} \leq \bar{K}_i$ , where  $\bar{K}_i = (1 + \|y_0^i\|) \exp(K_i T) - 1$ .

The aggregative optimal control problem of interest is defined as follows:

$$(5.2) \quad \inf_{(z_i, u_i)_{i=1}^N \in \prod_{i=1}^N \mathcal{X}_i} \int_0^T \ell \left( \frac{1}{N} \sum_{i=1}^N \phi_i(z_i(t), u_i(t)) \right) dt.$$

It is a special case of problem (P) with  $m = 1$ ,  $\mathcal{E}_1 = \mathcal{E} = L^2(0, T; \mathbb{R}^k)$ , and

$$\begin{aligned} g_i: \quad (z_i, u_i) \in \mathcal{X}_i &\mapsto (t \in (0, T) \mapsto \phi_i(z_i(t), u_i(t))) \in L^2(0, T; \mathbb{R}^k) \\ f: \quad y \in L^2(0, T; \mathbb{R}^k) &\mapsto \int_0^T \ell(y(t)) dt. \end{aligned}$$

Problem (5.2) can be seen as a nonconvex optimal control problem with state variable  $(z_i)_{i=1}^N$ . It finds application in energy management, in the situations mentioned in the introduction and in particular those involving storage devices, for which the dynamics of the state-of-charge must be taken into account. Once again we refer the reader to [41], which considers a convex stochastic aggregative optimal control problem. In general, only dynamic-programming-based methods can provide global solutions to nonlinear optimal control problems. They are not applicable here because of the high dimension of the state variable, equal to  $Nn$ .

It is easy to verify that  $\nabla f$  is continuously differentiable and that  $f$  and  $\nabla f$  are Lipschitz-continuous with moduli  $\sqrt{T}L_\ell$  and  $L_{\nabla\ell}$ , respectively. Let  $\hat{K}_i$  be an upper bound of  $\sup_{u_i \in U_i} \|u_i\|$ , for all  $i \in 1, \dots, N$ . Then  $g_i(\mathcal{X}_i)$  is bounded in  $L^2(0, T; \mathbb{R}^k)$ , with diameter bounded by  $2\sqrt{T}R_i(\hat{K}_i + \hat{K}_i)$ . Therefore, Assumption A is satisfied. If  $\ell$  is convex, then  $f$  is also convex and then Assumption B holds true. Let us verify Assumption C. Given  $y \in G(\mathcal{X})$ , the problem (3.1) to be solved at each iteration of the SFW algorithm reads

$$(5.3) \quad \inf_{(z_i, u_i) \in \mathcal{X}_i} \int_0^T \langle \nabla \ell(y(t)), \phi_i(z_i(t), u_i(t)) \rangle dt.$$

This is an optimal control with state variable  $z_i$ , which falls into the class of problems introduced in [21, Chapter III, Theorem 4.1] and therefore possesses a solution. If the dimension of the state variable,  $n$ , is small, then it can be solved by dynamic programming. We refer the reader to [19].

**5.4. Supervised learning problems.** We describe and discuss here two applications of problem (P) in the context of supervised learning.

*Neural networks with one hidden layer.* We refer the reader to [11, 37, 36]. Consider a neural network of the form  $\frac{1}{N} \sum_{i=1}^N \sigma_*(\mathbf{a}, x_i)$ , where  $\mathbf{a} \in \mathbb{R}^d$  is the feature vector,  $x = (x_i)_{i=1}^N \in (\mathbb{R}^D)^N$  are the network parameters (to be optimized), and  $\sigma_*: \mathbb{R}^d \times \mathbb{R}^D \rightarrow \mathbb{R}$  an activation function. We consider a loss function  $\varphi: \mathbb{R} \rightarrow \mathbb{R}_+$ . Given a data set  $(\mathbf{a}_j, b_j)_{j=1}^M \in (\mathbb{R}^d \times \mathbb{R})^M$ , the learning problem of interest writes

$$(5.4) \quad \inf_{(x_i)_{i=1}^N \in (\mathbb{R}^D)^N} \frac{1}{M} \sum_{j=1}^M \varphi \left( b_j - \frac{1}{N} \sum_{i=1}^N \sigma_*(\mathbf{a}_j, x_i) \right).$$

It is of the form (P), with  $\mathcal{E} = \mathbb{R}^M$ ,  $\mathcal{E}_j = \mathbb{R}$ ,  $f_j(y_j) = \varphi(b_j - y_j)/M$ ,  $g_{ij}(x_i) = \sigma_*(\mathbf{a}_j, x_i)$ . Assume that the set  $\{\sigma_*(\mathbf{a}_j, x) \mid x \in \mathbb{R}^D, j \in \{1, \dots, M\}\}$  has a bounded diameter  $\bar{d}$ . Assume moreover that  $\varphi$  is continuously differentiable and that  $\nabla\varphi$  is  $L_{\nabla\varphi}$ -Lipschitz continuous. Then Assumption A is satisfied and we have  $D_i = L_{\nabla\varphi}\bar{d}^2$ ,

for the coefficients  $D_i$  introduced in (4.1). Therefore, by Theorem 4.4, the optimality gap is bounded by

$$\frac{(M \wedge N)L_{\nabla\varphi}\bar{d}^2}{2N^2}.$$

Moreover, if  $\varphi$  is convex, then Assumption B holds true. The resolution of the sub-problems (3.1) is not easy in general, we refer the reader to [14] where the linearized problems are shown to be solvable by second-order cone programming in the case of ReLu activation functions.

Note that we are here in the symmetric case, as defined at the end of Section 2.2. The mean-field relaxation proposed in Lemma 2.8 was also utilized in [36, 37] for learning problems of the form (5.4). A gap estimate of order  $\mathcal{O}(1/N)$  is demonstrated, in the case of a quadratic loss function  $\varphi$ , see [37, Prop. 1]. Our gap estimate is more general see  $\nabla\varphi$  is only supposed to be Lipschitz; moreover, it is more precise in the case of an overparametrized network (i.e. when  $M < N$ ), since then it is of order  $\mathcal{O}(M/N^2)$ .

*Sparse reconstruction.* Another important learning example is the sparse reconstruction with the  $\ell_0$ -penalty, see [34, 33]. Let  $D$  be a  $M$  by  $N$  dictionary matrix. The objective is to approximate the observed vector  $x \in \mathbb{R}^M$  by a sparse linear combination of the columns of  $D$ . Following [33, Eq. 5.6], we are interested in the following least square problem with the  $\ell_0$ -penalty:

$$\inf_{\alpha \in \mathbb{R}^N} \frac{1}{2} \|x - D\alpha\|^2 + \beta \|\alpha\|_{\ell_0} = \inf_{\alpha \in \mathbb{R}^N} \frac{1}{2} \sum_{j=1}^M \left( x_j - \sum_{i=1}^N D_{ji} \alpha_i \right)^2 + \beta \sum_{i=1}^N \mathbf{1}_{\mathbb{R} \setminus \{0\}}(\alpha_i),$$

where  $\beta$  is a constant and  $\|\alpha\|_{\ell_0}$  counts the number of non-zero entries in a vector  $\alpha$ . Adding constraints of the form  $\alpha_i \in [u_i, v_i]$  to the problem, it is easy to see that Assumptions A and B are satisfied. The subproblems (3.1) are here of the form

$$\inf_{\alpha_i \in [u_i, v_i]} z\alpha_i + \mathbf{1}_{\mathbb{R} \setminus \{0\}}(\alpha_i)$$

for some real number  $z$ . One can show that the solution necessarily lies in  $\{u_i, v_i, 0\}$ , thus it is easy to compute.

Finally, let us mention other applications of the problem (P) in convex framework, for instance, the “sharing problem” in [9], Lasso regression in [20] and the dual problem of a linear support vector machine (SVM) in [42, 20].

**6. Numerical test.** In this section we provide numerical results for a mixed-integer linear quadratic problem of the form (P). Let  $A$  be a real  $M \times N$  matrix and let  $\bar{y} \in \mathbb{R}^M$ . Consider the following problem:

$$\text{(MIQP)} \quad \min_{x \in \{0,1\}^N} J(x) := \frac{1}{N^2} \|Ax - \bar{y}\|_{\mathbb{R}^M}^2 = \sum_{j=1}^M \left( \frac{1}{N} \sum_{i=1}^N A_{ji} x_i - \frac{\bar{y}_j}{N} \right)^2.$$

Problem (MIQP) has the form (P), with  $f_j(y_j) = (y_j - \frac{\bar{y}_j}{N})^2$  for  $1 \leq j \leq M$ , and  $g_{ij}(x_i) = A_{ji} x_i$  for  $1 \leq i \leq N$ ,  $1 \leq j \leq M$ . Moreover, Assumption A is satisfied with  $\tilde{L}_j = 2$  and  $d_{ij} = |A_{ji}|$ . Thus  $C_1 = \frac{2}{N} \sum_{i=1}^N \sum_{j=1}^M |A_{ji}|$ . Due to the linearity of  $g_{ij}$ , the randomized problem coincides with the minimization problem of  $J$  on  $[0, 1]^N$ , which

is a convex linear-quadratic program that can be solved with independent methods; thus it is easy here to obtain a precise estimate of  $\mathcal{J}^*$ .

In the numerical simulation, we draw the parameters  $A_{ji}$  according to the uniform distribution on the interval  $[0, 1]$  while  $y_j$  is drawn according to the uniform distribution on  $[0, N/2]$ . Thus,  $C_1 \approx M$  and the gap estimate is given by  $\frac{C_1}{2N} \approx 0.5$ . We perform our numerical experiments on a laptop with one Intel Core i5-8250U processor (4 cores) at 1.60 GHz and 8 GB RAM.

The first experiment is a comparison of Algorithm 2 with an open source solver, SCIP, [7] and a commercial solver, GUROBI, [23]. As mentioned before, the dual (randomized) problem is a convex linear-quadratic program. We can compute  $\mathcal{J}^*$  easily by solver GUROBI. Table 1 shows the value  $\mathcal{J}^*$  and results of (MIQP) obtained from SCIP, GUROBI and Algorithm 2, for different values of  $M, N$  ranging from 100 to 3200. In Table 1, “Nan” indicates that the solver has failed to return a result or that computation time has exceeded one hour. Denote by  $v_s$  the result of Algorithm 2. The indicated gap is a relative gap, in percent, defined by  $(v_s - \mathcal{J}^*)/\mathcal{J}^*$ . We can observe that the relative gap decreases as  $N$  increases, which is consistent with the randomized gap (2.2). The last three columns of Table 1 show that Algorithm 2 is competitive in terms of execution time, in comparison with SCIP and GUROBI. Finally, observe that for  $N = M = 3200$ , none of the two solvers could solve the problems while Algorithm 2 has provided a solutions in approximately 6 minutes.

$N = M$	$\mathcal{J}^*$	SCIP	GUR.	SFW		SCIP	GUR.	SFW
		value	value	value	gap in %	time in seconds		
100	2.077	2.077	2.077	2.136	2.870	0.88	0.20	0.03
200	4.120	4.120	4.120	4.159	0.956	5.99	0.69	0.09
400	7.871	7.871	7.871	7.904	0.430	87.78	7.90	0.91
800	15.953	Nan	15.954	15.966	0.079	Nan	10.63	6.18
1600	32.045	Nan	32.048	32.0585	0.042	Nan	81.41	42.51
3200	64.717	Nan	Nan	64.724	0.012	Nan	Nan	330.95

Table 1: Comparison of the approximate values and execution times obtained with SCIP, GUROBI and Algorithm 2 for problem (MIQP) with  $M = N = 100, 200, 400, 800, 1600$  and  $3200$ . In Algorithm 2, we take  $n_k = 1$  and  $K = 2N$  iterations.

The second experiment is on the basic Frank-Wolfe algorithm 1 and its stochastic version 2. In this experiment, we fix  $M = N = 1000$ . Figure 1 shows the outcome of the basic Frank-Wolfe algorithm 1 with 200 iterations. The left sub-figure shows the evolution of  $\gamma_k$  for  $\omega_k = 2/(k+2)$  (green curve) and for  $\omega_k$  determined by line search (3.6) (red curve). A sub-linear rate of convergence is observed (note that logarithmic scales are employed for both axes). The right sub-figure represents the evolution of  $J(X^k) - \mathcal{J}^*$ , where  $X^k$  is a random variable with distribution  $\mu^k$ . For both choices of  $\omega_k$ , approximate solutions to the problems are simulated, with a gap smaller than  $10^{-3}$ , significantly smaller than the gap estimate  $\frac{C_1}{2N}$ . The line search approach is quicker than the approach with  $\omega_k = \frac{2}{k+2}$ .

Figure 2 shows the outcome of Algorithm 2 (with the modification suggested in Remark 3.9), for different (constant) choices of  $n_k$  with 200 iterations, for two different stepsize rules ( $\omega_k = 2/(k+2)$  on the left, line search on the right). Since the algorithm is stochastic, we have tested it 50 times to evaluate its efficiency; the curves represent

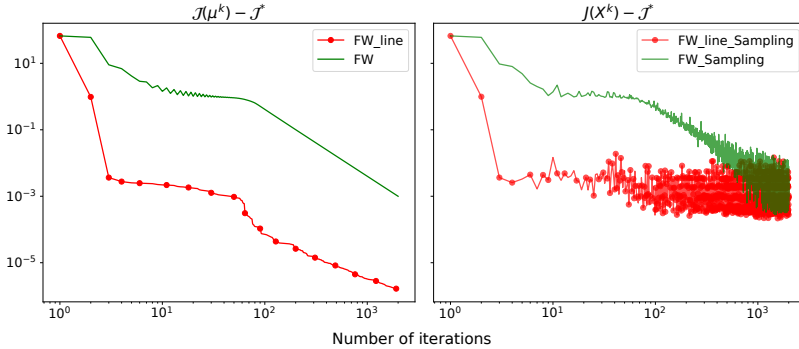


Fig. 1: MIQP by Algorithm 1, 2000 iterations, with  $\omega_k = 2/(k + 2)$  and line search (3.6).

the average value of  $\gamma_k$ . The standard deviation (for these 8 instances of the SFW method) is displayed on Figure 3. In all cases, an average value of the gap significantly smaller than  $\frac{C_1}{2N}$  can be reached; the standard deviation is also significantly smaller than  $\frac{C_1}{2N}$  at the last iterations. There is a benefit (both in expectation and standard deviation) in increasing the number of simulations  $n_k$  (note that the choice  $n_k = 1000$  is much smaller the rule suggested by Corollary 3.8). Yet the convergence is slower in comparison with the basic Frank-Wolfe algorithm, which can be explained by the use of the selection method at each iteration.

**7. Conclusion.** We have investigated a large-scale and aggregative optimization problem and its relaxation. New error bounds for the relaxation gap have been obtained. We have proposed a tractable algorithm for its resolution with a detailed convergence analysis relying on concentration inequalities. Assuming that an efficient method for the resolution of the subproblems is available, the implementation of our stochastic Frank-Wolfe method is easy.

Future research will focus on refinements of the selection method, allowing the computation of  $\mathcal{O}(q \wedge N/N^2)$ -solutions. We also aim at working on more complex problems, involving for example convex constraints on the aggregate, as for example the resource allocation problems mentioned in the introduction. Such constraints could be handled with extensions of the Frank-Wolfe algorithm for non-smooth costs as those proposed in [44, 50]. Finally, we intend to apply our method to large-scale optimal control problems, such as nonconvex variants of the problem investigated in [41].

#### Appendix A. Concentration inequalities and other technical lemmas.

PROPOSITION A.1. Consider  $T$  real-valued random variables  $(Y_t)_{t=1,\dots,T}$ . Let  $(\mathcal{F}_t)_{t=1,\dots,T}$  denote the associated filtration ( $\mathcal{F}_0$  is the trivial  $\sigma$ -algebra). Let  $Z_t = \mathbb{E}[Y_t^2 | \mathcal{F}_{t-1}]$  and let  $S_T = \sum_{t=1}^T Y_t$ . Assume the following:

$$(A.1) \quad (i) \quad \mathbb{E}[Y_t | \mathcal{F}_{t-1}] = 0, \quad (ii) \quad Y_t \leq m, \quad (iii) \quad \sum_{t'=1}^T Z_{t'} \leq v, \quad a.s.$$

for all  $t = 1, \dots, T$  and for some constants  $m$  and  $v$ . Then,  $\mathbb{E}[S_T^2] \leq v$ . Moreover, for

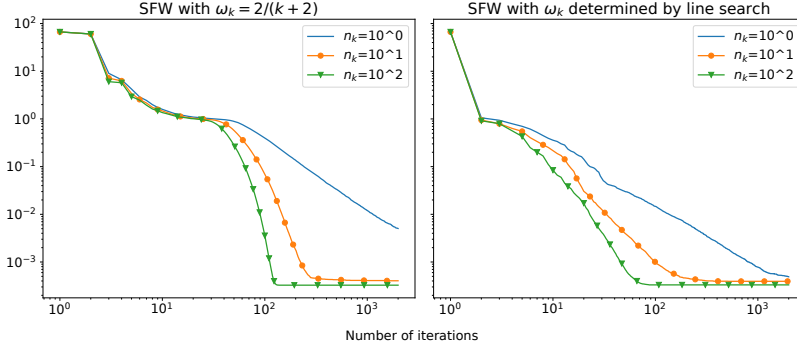


Fig. 2: MIQP by Algorithm 2 with 2000 iterations, expectation of the gap.

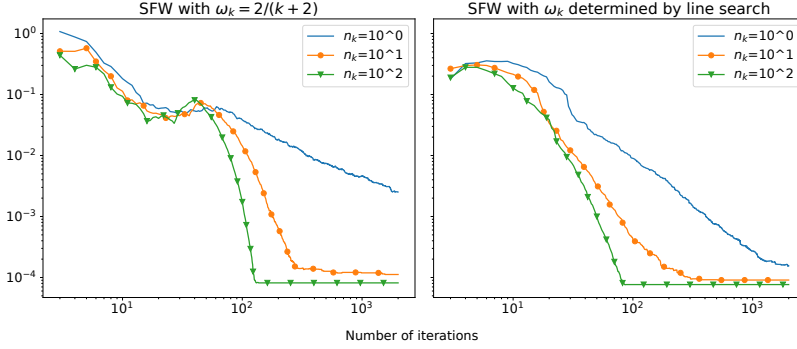


Fig. 3: MIQP by Algorithm 2 with 2000 iterations, standard deviation of the gap.

any  $\epsilon > 0$ ,

$$(A.2) \quad \mathbb{P}[S_T \geq \epsilon] \leq \exp\left(-\frac{\epsilon^2}{2(v + \epsilon m/3)}\right).$$

*Proof.* The estimate of  $\mathbb{E}[S_T^2]$  can be easily obtained by induction. For the estimate of  $\mathbb{P}[S_T \geq \epsilon]$ , see [15, Theorem 7].  $\square$

As a corollary, we obtain the following McDiarmid's inequality of "variance type".

**COROLLARY A.2.** *Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space and let  $(\Omega_i)_{i=1, \dots, n}$  be  $n$  measurable subsets of  $\Omega$ . Let  $X = (X_i)_{i=1, \dots, n}$  be  $n$  independent random variables valued respectively in  $(\Omega_i)_{i=1, \dots, n}$ . Consider a measurable function  $f: \prod_{i=1}^n \Omega_i \rightarrow \mathbb{R}$  and real constants  $m \in \mathbb{R}$  and  $(v_i)_{i=1, \dots, n}$  such that*

$$\text{Var}[f(X_i, x_{-i})] \leq v_i^2, \quad a.s., \quad |f(X_i, x_{-i}) - \mathbb{E}[f(X_i, x_{-i})]| \leq m, \quad a.s.,$$

for all  $i = 1, \dots, n$  and for all  $x_{-i} \in \left(\prod_{j=1}^{i-1} \Omega_j\right) \times \left(\prod_{j=i+1}^n \Omega_j\right)$ . Then, for any  $\epsilon > 0$ ,

$$(A.3) \quad \mathbb{P}\left[f(\mathbf{x}) - \mathbb{E}[f(\mathbf{x})] \geq \epsilon\right] \leq \exp\left(-\frac{\epsilon^2}{2\left(\sum_{i=1}^n v_i^2 + \frac{m\epsilon}{3}\right)}\right).$$

*Proof.* Define  $Y_t = \mathbb{E}[f(X) \mid X_1, \dots, X_t] - \mathbb{E}[f(X) \mid X_1, \dots, X_{t-1}]$  and apply Proposition A.1.  $\square$

LEMMA A.3. For all  $k \in \mathbb{N}$ , denote  $\omega_k = \frac{2}{k+2}$ . Let  $(u_k)_{k \in \mathbb{N}}$  and  $(\gamma_k)_{k \in \mathbb{N}}$  be two sequences of real numbers. Assume that there exists a positive number  $C$  such that

$$(A.4) \quad \gamma_{k+1} \leq (1 - \omega_k)\gamma_k + C\omega_k^2 + u_k,$$

for all  $k \in \mathbb{N}$ . Then, for all  $K \in \mathbb{N}^*$ ,

$$(A.5) \quad \gamma_K \leq \frac{4C}{K} + \sum_{k=0}^{K-1} \frac{(k+1)(k+2)}{K(K+1)} u_k.$$

*Proof.* We prove this lemma by induction on  $K$ . We have  $\omega_0 = 1$ , thus taking  $k = 0$  in (A.4), we obtain that  $\gamma_1 \leq C + u_0$ , which proves the claim for  $K = 1$ . Let us assume that the claim holds true for some  $K \in \mathbb{N}^*$ . We deduce from (A.4) that

$$\begin{aligned} \gamma_{K+1} &\leq \left( \frac{1}{K+2} + \frac{1}{(K+2)^2} \right) 4C + \frac{K}{K+2} \left( \sum_{k=0}^{K-1} \frac{(k+1)(k+2)}{K(K+1)} u_k \right) + u_K \\ &\leq \frac{4C}{K+1} + \sum_{k=0}^K \frac{(k+1)(k+2)}{(K+1)(K+2)} u_k. \end{aligned}$$

Therefore the claim holds for  $K+1$ . This concludes the proof.  $\square$

LEMMA A.4. Let  $A$ ,  $B$ , and  $C$  be three random variables. Assume that  $B$  is independent of  $(A, C)$  and that  $B \sim \text{Bern}(\omega)$  for some  $\omega \in [0, 1]$ . Let  $F$  be a real-valued function of  $(A, B, C)$ . Assume that  $|F(A, 1, C) - F(A, 0, C)| \leq \delta$ , a.s. Finally, define  $U = \mathbb{E}[F(A, B, C) \mid A, B] - \mathbb{E}[F(A, B, C) \mid A]$ . Then,

$$\mathbb{E}[U \mid A] = 0, \quad U \leq \delta, \quad \mathbb{E}[U^2 \mid A] \leq \omega(1 - \omega)\delta^2, \quad \text{a.s.}$$

*Proof.* The equality  $\mathbb{E}[U \mid A] = 0$  is trivial. We have  $U = \mathbb{E}[Z \mid A, B]$ , where

$$Z = F(A, B, C) - \mathbb{E}[F(A, B, C) \mid A, C].$$

It is easy to verify that  $Z \leq \delta$ , a.s., which implies that  $\mathbb{E}[U \mid A] = \mathbb{E}[Z \mid A] \leq \delta$ . The first inequality is proved. For the second inequality, we first note that

$$\mathbb{E}[Z^2 \mid A, C] = \omega(1 - \omega)(F(A, 1, C) - F(A, 0, C))^2,$$

as can be easily verified. Thus  $\mathbb{E}[Z \mid A] \leq \omega(1 - \omega)\delta^2$ . Next by Jensen's inequality, we have  $U^2 \leq \mathbb{E}[Z^2 \mid A, B]$ . Therefore,

$$\mathbb{E}[U^2 \mid A] \leq \mathbb{E}[\mathbb{E}[Z^2 \mid A, B] \mid A] = \mathbb{E}[Z^2 \mid A] \leq \omega(1 - \omega)\delta^2,$$

as was to be proved.  $\square$

The following lemma is an elementary property of the conditional expectation. For the sake of simplicity, we only state it (and prove it) with discrete random variables.

LEMMA A.5. Let  $X$ ,  $Y$ , and  $Z$  be three random variables. Assume that  $Y$  and  $Z$  are discrete and that  $Z$  is independent of  $(X, Y)$ . Then,  $\mathbb{E}[X \mid Y, Z] = \mathbb{E}[X \mid Y]$ .



*Proof.* By definition,  $\mathbb{E}[X | Y, Z] = \phi(Y, Z)$ , where  $\phi$  is defined as follows: for any pair  $(y, z)$  such that  $\mathbb{P}[Y = y \text{ and } Z = z] \neq 0$ ,

$$\phi(y, z) = \frac{\mathbb{E}[X \mathbf{1}_{Y=y} \mathbf{1}_{Z=z}]}{\mathbb{P}[Y = y \text{ and } Z = z]} = \frac{\mathbb{E}[X \mathbf{1}_{Y=y}]}{\mathbb{P}[Y = y]},$$

since  $Z$  is independent of  $(X, Y)$ . Thus  $\phi$  does not depend on  $Z$  and the result follows.  $\square$

#### REFERENCES

- [1] M. Atallah, M. Bykova, J. Li, K. Frikken, and M. Topkara. Private collaborative forecasting and benchmarking. In *Proceedings of the 2004 ACM workshop on Privacy in the electronic society*, pages 103–114, 2004.
- [2] J.P. Aubin and I. Ekeland. Estimates of the duality gap in nonconvex optimization. *Mathematics of Operations Research*, 1(3):225–245, 1976.
- [3] H.H. Bauschke and P.L. Combettes. *Convex analysis and monotone operator theory in Hilbert spaces*, volume 408. Springer, 2011.
- [4] O. Beaude, P. Benchimol, S. Gaubert, P. Jacquot, and N. Oudjane. A privacy-preserving method to optimize distributed resource allocation. *SIAM Journal on Optimization*, 30(3):2303–2336, 2020.
- [5] A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- [6] A. Beck and L. Tetruashvili. On the convergence of block coordinate descent type methods. *SIAM Journal on Optimization*, 23(4):2037–2060, 2013.
- [7] K. Bestuzheva, M. Besançon, W.-K. Chen, A. Chmiela, T. Donkiewicz, J. van Doornmalen, L. Eifler, O. Gaul, G. Gamrath, A. Gleixner, et al. The scip optimization suite 8.0. *arXiv preprint arXiv:2112.08872*, 2021.
- [8] P. Bonami, M. Kilinç, and J. Linderoth. Algorithms and software for convex mixed integer nonlinear programs. In *Mixed integer nonlinear programming*, pages 1–39. Springer, 2012.
- [9] S. Boyd, N. Parikh, and E. Chu. *Distributed optimization and statistical learning via the alternating direction method of multipliers*. Now Publishers Inc, 2011.
- [10] J. Cassels. Measures of the non-convexity of sets and the Shapley–Folkman–Starr theorem. In *Mathematical Proceedings of the Cambridge Philosophical Society*, volume 78, pages 433–436. Cambridge University Press, 1975.
- [11] L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [12] C. Coey, M. Lubin, and J.P. Vielma. Outer approximation with conic certificates for mixed-integer convex problems. *Mathematical Programming Computation*, 12(2):249–293, 2020.
- [13] L. Condat. Fast projection onto the simplex and the  $l_1$  ball. *Mathematical Programming*, 158(1):575–585, 2016.
- [14] A. d’Aspremont and M. Pilanci. Global convergence of frank wolfe on one hidden layer networks. *arXiv preprint arXiv:2002.02208*, 2020.
- [15] B. Delyon. Exponential inequalities for dependent processes. Technical report, October 2015.
- [16] L. Ding and M. Udell. Frank-Wolfe style algorithms for large scale optimization. In *Large-Scale and Distributed Optimization*, pages 215–245. Springer, 2018.
- [17] J.C. Dunn and S. Harshbarger. Conditional gradient algorithms with open loop step size rules. *Journal of Mathematical Analysis and Applications*, 62(2):432–444, 1978.
- [18] J. Fadili, C. Molinari, and A. Silveti-Falls. Inexact and stochastic generalized conditional gradient with augmented lagrangian and proximal step. *Journal of Nonsmooth Analysis and Optimization*, 2, 2021.
- [19] M. Falcone and R. Ferretti. *Semi-Lagrangian approximation schemes for linear and Hamilton—Jacobi equations*. SIAM, 2013.
- [20] O. Fercoq and P. Richtárik. Optimization in high dimensions via accelerated, parallel, and proximal coordinate descent. *Siam review*, 58(4):739–771, 2016.
- [21] W.H. Fleming and R.W. Rishel. *Deterministic and stochastic optimal control*, volume 1 of *Appl. Math. (N. Y.)*. Springer, New York, 1975.
- [22] D. Ghaderyan, F.L. Pereira, and A.P. Aguiar. A fully distributed method for distributed multiagent system in a microgrid. *Energy Reports*, 7:2294–2301, 2021.
- [23] LLC Gurobi Optimization. Gurobi optimizer reference manual, 2018.

- [24] H. Hassani, A. Karbasi, A. Mokhtari, and Z. Shen. Stochastic conditional gradient++:(non) convex minimization and continuous submodular maximization. *SIAM Journal on Optimization*, 30(4):3315–3344, 2020.
- [25] E. Hazan and H. Luo. Variance-reduced and projection-free stochastic optimization. In *International Conference on Machine Learning*, pages 1263–1271. PMLR, 2016.
- [26] J.B. Hiriart-Urruty and C. Lemaréchal. *Convex analysis and minimization algorithms I: Fundamentals*, volume 305. Springer science & business media, 1996.
- [27] P. Jacquot, O. Beaude, S. Gaubert, and N. Oudjane. Analysis and implementation of an hourly billing mechanism for demand response management. *IEEE Transactions on Smart Grid*, 10(4):4265–4278, 2018.
- [28] M. Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *International Conference on Machine Learning*, pages 427–435. PMLR, 2013.
- [29] T. Kerdreux, I. Colin, and A. D’Aspremont. Stable bounds on the duality gap of separable nonconvex optimization problems. *Mathematics of Operations Research*, 2022.
- [30] W. Krichene, A. Bayen, and P.L. Bartlett. Accelerated mirror descent in continuous and discrete time. *Advances in neural information processing systems*, 28, 2015.
- [31] S. Lacoste-Julien and M. Jaggi. On the global linear convergence of frank-wolfe optimization variants. *Advances in neural information processing systems*, 28, 2015.
- [32] F. Locatello, A. Yurtsever, O. Fercoq, and V. Cevher. Stochastic Frank-Wolfe for composite convex minimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [33] J. Mairal, F. Bach, and J. Ponce. Sparse modeling for image and vision processing. *arXiv preprint arXiv:1411.3230*, 2014.
- [34] S. Mallat. *A wavelet tour of signal processing*. Elsevier, 2009.
- [35] C. McDiarmid. On the method of bounded differences. *Surveys in combinatorics*, 141(1):148–188, 1989.
- [36] S. Mei, T. Misiakiewicz, and A. Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In *Conference on Learning Theory*, pages 2388–2464. PMLR, 2019.
- [37] S. Mei, A. Montanari, and P.-M. Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- [38] A. Mokhtari, H. Hassani, and A. Karbasi. Stochastic conditional gradient methods: From convex minimization to submodular maximization. *Journal of machine learning research*, 2020.
- [39] F. Pacaud. *Decentralized optimization for energy efficiency under stochasticity*. PhD Thesis, Université Paris-Est, October 2018.
- [40] R.T. Rockafellar. *Convex analysis*. Princeton university press, 2015.
- [41] A. Seguret, C. Alasseur, J.F. Bonnans, A. De Paola, N. Oudjane, and V. Trovato. Decomposition of high dimensional aggregative stochastic control problems. *arXiv preprint arXiv:2008.09827*, 2020.
- [42] S. Shalev-Shwartz and A. Tewari. Stochastic methods for  $l_1$  regularized loss minimization. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 929–936, 2009.
- [43] R.-H. Shi, Y. Mu, H. Zhong, J. Cui, and S. Zhang. Secure multiparty quantum computation for summation and multiplication. *Scientific reports*, 6(1):1–9, 2016.
- [44] A. Silveti-Falls, C. Molinari, and J. Fadili. Generalized conditional gradient with augmented Lagrangian for composite minimization. *SIAM Journal on Optimization*, 30(4):2687–2725, 2020.
- [45] R.M. Starr. Quasi-equilibria in markets with non-convex preferences. *Econometrica: journal of the Econometric Society*, pages 25–38, 1969.
- [46] T. Tang, K. Balasubramanian, and T.C.M. Lee. High-probability bounds for robust stochastic frank-wolfe algorithm. In *The 38th Conference on Uncertainty in Artificial Intelligence*, 2022.
- [47] A.B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, 2009.
- [48] C. Villani. *Topics in optimal transportation*, volume 58. American Mathematical Soc., 2003.
- [49] M. Wang. Vanishing price of decentralization in large coordinative nonconvex optimization. *SIAM Journal on Optimization*, 27(3):1977–2009, 2017.
- [50] A. Yurtsever, O. Fercoq, and V. Cevher. A conditional-gradient-based augmented lagrangian framework. In *International Conference on Machine Learning*, pages 7272–7281, 2019.
- [51] A. Yurtsever, S. Sra, and V. Cevher. Conditional gradient methods via stochastic path-integrated differential estimator. In *International Conference on Machine Learning*, pages 7282–7291. PMLR, 2019.