



**HAL**  
open science

# Stochastic Majorization-Minimization Algorithms for Large-Scale Optimization

Julien Mairal

► **To cite this version:**

Julien Mairal. Stochastic Majorization-Minimization Algorithms for Large-Scale Optimization. NIPS 2013 - Advances in Neural Information Processing Systems, Dec 2013, South Lake Tahoe, United States. pp.2283-2291. hal-00835840v2

**HAL Id: hal-00835840**

**<https://inria.hal.science/hal-00835840v2>**

Submitted on 10 Sep 2013

**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.

---

# Stochastic Majorization-Minimization Algorithms for Large-Scale Optimization

---

**Julien Mairal**

LEAR Project-Team - INRIA Grenoble  
julien.mairal@inria.fr

## Abstract

Majorization-minimization algorithms consist of iteratively minimizing a majorizing surrogate of an objective function. Because of its simplicity and its wide applicability, this principle has been very popular in statistics and in signal processing. In this paper, we intend to make this principle scalable. We introduce a stochastic majorization-minimization scheme which is able to deal with large-scale or possibly infinite data sets. When applied to convex optimization problems under suitable assumptions, we show that it achieves an expected convergence rate of  $O(1/\sqrt{n})$  after  $n$  iterations, and of  $O(1/n)$  for strongly convex functions. Equally important, our scheme almost surely converges to stationary points for a large class of non-convex problems. We develop several efficient algorithms based on our framework. First, we propose a new stochastic proximal gradient method, which experimentally matches state-of-the-art solvers for large-scale  $\ell_1$ -logistic regression. Second, we develop an online DC programming algorithm for non-convex sparse estimation. Finally, we demonstrate the effectiveness of our approach for solving large-scale structured matrix factorization problems.

## 1 Introduction

Majorization-minimization [15] is a simple optimization principle for minimizing an objective function. It consists of iteratively minimizing a surrogate that upper-bounds the objective, thus monotonically driving the objective function value downhill. This idea is used in many existing procedures. For instance, the expectation-maximization (EM) algorithm (see [5, 21]) builds a surrogate for a likelihood model by using Jensen’s inequality. Other approaches can also be interpreted under the majorization-minimization point of view, such as DC programming [8], where “DC” stands for difference of convex functions, variational Bayes techniques [28], or proximal algorithms [1, 23, 29].

In this paper, we propose a stochastic majorization-minimization algorithm, which is suitable for solving large-scale problems arising in machine learning and signal processing. More precisely, we address the minimization of an expected cost—that is, an objective function that can be represented by an expectation over a data distribution. For such objectives, online techniques based on stochastic approximations have proven to be particularly efficient, and have drawn a lot of attraction in machine learning, statistics, and optimization [3–6, 9–12, 14, 16, 17, 19, 22, 24–26, 30].

Our scheme follows this line of research. It consists of iteratively building a surrogate of the expected cost when only a single data point is observed at each iteration; this data point is used to update the surrogate, which in turn is minimized to obtain a new estimate. Some previous works are closely related to this scheme: the online EM algorithm for latent data models [5, 21] and the online matrix factorization technique of [19] involve for instance surrogate functions updated in a similar fashion. Compared to these two approaches, our method is targeted to more general optimization problems.

Another related work is the incremental majorization-minimization algorithm of [18] for finite training sets; it was indeed shown to be efficient for solving machine learning problems where storing

dense information about the past iterates can be afforded. Concretely, this incremental scheme requires to store  $O(pn)$  values, where  $p$  is the variable size, and  $n$  is the size of the training set.<sup>1</sup> This issue was the main motivation for us for proposing a stochastic scheme with a memory load independent of  $n$ , thus allowing us to possibly deal with infinite data sets, or a huge variable size  $p$ .

We study the convergence properties of our algorithm when the surrogates are strongly convex and chosen among the class of *first-order surrogate functions* introduced in [18], which consist of approximating the possibly non-smooth objective up to a smooth error. When the objective is convex, we obtain expected convergence rates that are asymptotically optimal, or close to optimal [14, 22]. More precisely, the convergence rate is of order  $O(1/\sqrt{n})$  in a finite horizon setting, and  $O(1/n)$  for a strongly convex objective in an infinite horizon setting. Our second analysis shows that for *non-convex* problems, our method almost surely converges to a set of stationary points under suitable assumptions. We believe that this result is equally valuable as convergence rates for convex optimization. To the best of our knowledge, the literature on stochastic non-convex optimization is rather scarce, and we are only aware of convergence results in more restricted settings than ours—see for instance [3] for the stochastic gradient descent algorithm, [5] for online EM, [19] for online matrix factorization, or [9], which provides stronger guarantees, but for unconstrained smooth problems.

We develop several efficient algorithms based on our framework. The first one is a new stochastic proximal gradient method for composite or constrained optimization. This algorithm is related to a long series of work in the convex optimization literature [6, 10, 12, 14, 16, 22, 25, 30], and we demonstrate that it performs as well as state-of-the-art solvers for large-scale  $\ell_1$ -logistic regression [7]. The second one is an online DC programming technique, which we demonstrate to be better than batch alternatives for large-scale non-convex sparse estimation [8]. Finally, we show that our scheme can address efficiently structured sparse matrix factorization problems in an online fashion, and offers new possibilities to [13, 19] such as the use of various loss or regularization functions.

This paper is organized as follows: Section 2 introduces first-order surrogate functions for batch optimization; Section 3 is devoted to our stochastic approach and its convergence analysis; Section 4 presents several applications and numerical experiments, and Section 5 concludes the paper.

## 2 Optimization with First-Order Surrogate Functions

Throughout the paper, we are interested in the minimization of a continuous function  $f : \mathbb{R}^p \rightarrow \mathbb{R}$ :

$$\min_{\theta \in \Theta} f(\theta), \quad (1)$$

where  $\Theta \subseteq \mathbb{R}^p$  is a convex set. The majorization-minimization principle consists of computing a majorizing surrogate  $g_n$  of  $f$  at iteration  $n$  and updating the current estimate by  $\theta_n \in \arg \min_{\theta \in \Theta} g_n(\theta)$ . The success of such a scheme depends on how well the surrogates approximate  $f$ . In this paper, we consider a particular class of surrogate functions introduced in [18] and defined as follows:

### Definition 2.1 (Strongly Convex First-Order Surrogate Functions).

Let  $\kappa$  be in  $\Theta$ . We denote by  $\mathcal{S}_{L,\rho}(f, \kappa)$  the set of  $\rho$ -strongly convex functions  $g$  such that  $g \geq f$ ,  $g(\kappa) = f(\kappa)$ , the approximation error  $g - f$  is differentiable, and the gradient  $\nabla(g - f)$  is  $L$ -Lipschitz continuous. We call the functions  $g$  in  $\mathcal{S}_{L,\rho}(f, \kappa)$  “first-order surrogate functions”.

Among the first-order surrogate functions presented in [18], we should mention the following ones:

- **Lipschitz Gradient Surrogates.**

When  $f$  is differentiable and  $\nabla f$  is  $L$ -Lipschitz,  $f$  admits the following surrogate  $g$  in  $\mathcal{S}_{2L,L}(f, \kappa)$ :

$$g : \theta \mapsto f(\kappa) + \nabla f(\kappa)^\top (\theta - \kappa) + \frac{L}{2} \|\theta - \kappa\|_2^2.$$

When  $f$  is convex,  $g$  is in  $\mathcal{S}_{L,L}(f, \kappa)$ , and when  $f$  is  $\mu$ -strongly convex,  $g$  is in  $\mathcal{S}_{L-\mu,L}(f, \kappa)$ . Minimizing  $g$  amounts to performing a classical gradient descent step  $\theta \leftarrow \kappa - \frac{1}{L} \nabla f(\kappa)$ .

- **Proximal Gradient Surrogates.**

Assume that  $f$  splits into  $f = f_1 + f_2$ , where  $f_1$  is differentiable,  $\nabla f_1$  is  $L$ -Lipschitz, and  $f_2$  is

<sup>1</sup>To alleviate this issue, it is possible to cut the dataset into  $\eta$  mini-batches, reducing the memory load to  $O(p\eta)$ , which remains cumbersome when  $p$  is very large.

convex. Then, the function  $g$  below is in  $\mathcal{S}_{2L,L}(f, \kappa)$ :

$$g : \theta \mapsto f_1(\kappa) + \nabla f_1(\kappa)^\top (\theta - \kappa) + \frac{L}{2} \|\theta - \kappa\|_2^2 + f_2(\theta).$$

When  $f_1$  is convex,  $g$  is in  $\mathcal{S}_{L,L}(f, \kappa)$ . If  $f_1$  is  $\mu$ -strongly convex,  $g$  is in  $\mathcal{S}_{L-\mu,L}(f, \kappa)$ . Minimizing  $g$  amounts to a proximal gradient step [1, 23, 29]:  $\theta \leftarrow \arg \min_{\theta} \frac{1}{2} \|\kappa - \frac{1}{L} \nabla f_1(\kappa) - \theta\|_2^2 + \frac{1}{L} f_2(\theta)$ .

• **DC Programming Surrogates.**

Assume that  $f = f_1 + f_2$ , where  $f_2$  is concave and differentiable,  $\nabla f_2$  is  $L_2$ -Lipschitz, and  $g_1$  is in  $\mathcal{S}_{L_1, \rho_1}(f_1, \kappa)$ . Then, the following function  $g$  is a surrogate in  $\mathcal{S}_{L_1+L_2, \rho_1}(f, \kappa)$ :

$$g : \theta \mapsto f_1(\theta) + f_2(\kappa) + \nabla f_2(\kappa)^\top (\theta - \kappa).$$

When  $f_1$  is convex,  $f_1 + f_2$  is a difference of convex functions, leading to a DC program [8].

With the definition of first-order surrogates and a basic “batch” algorithm in hand, we now introduce our main contribution: a stochastic scheme for solving large-scale problems.

### 3 Stochastic Optimization

As pointed out in [4], one is usually not interested in the minimization of an *empirical cost* on a finite training set, but instead in minimizing an *expected cost*. Thus, we assume from now on that  $f$  has the form of an expectation:

$$\min_{\theta \in \Theta} \left[ f(\theta) \triangleq \mathbb{E}_{\mathbf{x}}[\ell(\mathbf{x}, \theta)] \right], \quad (2)$$

where  $\mathbf{x}$  from some set  $\mathcal{X}$  represents a data point, which is drawn according to some unknown distribution, and  $\ell$  is a continuous loss function. As often done in the literature [22], we assume that the expectations are well defined and finite valued; we also assume that  $f$  is bounded below.

We present our approach for tackling (2) in Algorithm 1. At each iteration, we draw a training point  $\mathbf{x}_n$ , assuming that these points are i.i.d. samples from the data distribution. Note that in practice, since it is often difficult to obtain true i.i.d. samples, the points  $\mathbf{x}_n$  are computed by cycling on a randomly permuted training set [4]. Then, we choose a surrogate  $g_n$  for the function  $\theta \mapsto \ell(\mathbf{x}_n, \theta)$ , and we use it to update a function  $\bar{g}_n$  that behaves as an approximate surrogate for the expected cost  $f$ . The function  $\bar{g}_n$  is in fact a weighted average of previously computed surrogates, and involves a sequence of weights  $(w_n)_{n \geq 1}$  that will be discussed later. Then, we minimize  $\bar{g}_n$ , and obtain a new estimate  $\theta_n$ . For convex problems, we also propose to use averaging schemes, denoted by “option 2” and “option 3” in Alg. 1. Averaging is a classical technique for improving convergence rates in convex optimization [10, 22] for reasons that are clear in the convergence proofs.

---

**Algorithm 1** Stochastic Majorization-Minimization Scheme

---

**input**  $\theta_0 \in \Theta$  (initial estimate);  $N$  (number of iterations);  $(w_n)_{n \geq 1}$ , weights in  $(0, 1]$ ;

1: initialize the approximate surrogate:  $\bar{g}_0 : \theta \mapsto \frac{\rho}{2} \|\theta - \theta_0\|_2^2$ ;  $\bar{\theta}_0 = \theta_0$ ;  $\hat{\theta}_0 = \theta_0$ ;

2: **for**  $n = 1, \dots, N$  **do**

3: draw a training point  $\mathbf{x}_n$ ; define  $f_n : \theta \mapsto \ell(\mathbf{x}_n, \theta)$ ;

4: choose a surrogate function  $g_n$  in  $\mathcal{S}_{L, \rho}(f_n, \theta_{n-1})$ ;

5: update the approximate surrogate:  $\bar{g}_n = (1 - w_n) \bar{g}_{n-1} + w_n g_n$ ;

6: update the current estimate:

$$\theta_n \in \arg \min_{\theta \in \Theta} \bar{g}_n(\theta);$$

7: for option 2, update the averaged iterate:  $\hat{\theta}_n \triangleq (1 - w_{n+1}) \hat{\theta}_{n-1} + w_{n+1} \theta_n$ ;

8: for option 3, update the averaged iterate:  $\bar{\theta}_n \triangleq \frac{(1 - w_{n+1}) \bar{\theta}_{n-1} + w_{n+1} \theta_n}{\sum_{k=1}^{n+1} w_k}$ ;

9: **end for**

**output (option 1):**  $\theta_N$  (current estimate, no averaging);

**output (option 2):**  $\bar{\theta}_N$  (first averaging scheme);

**output (option 3):**  $\hat{\theta}_N$  (second averaging scheme).

---

We remark that Algorithm 1 is only practical when the functions  $\bar{g}_n$  can be parameterized with a small number of variables, and when they can be easily minimized over  $\Theta$ . Concrete examples are discussed in Section 4. Before that, we proceed with the convergence analysis.

### 3.1 Convergence Analysis - Convex Case

First, We study the case of convex functions  $f_n : \theta \mapsto \ell(\theta, \mathbf{x}_n)$ , and make the following assumption:

- (A) for all  $\theta$  in  $\Theta$ , the functions  $f_n$  are  $R$ -Lipschitz continuous. Note that for convex functions, this is equivalent to saying that subgradients of  $f_n$  are uniformly bounded by  $R$ .

Assumption (A) is classical in the stochastic optimization literature [22]. Our first result shows that with the averaging scheme corresponding to “option 2” in Alg. 1, we obtain an expected convergence rate that makes explicit the role of the weight sequence  $(w_n)_{n \geq 1}$ .

**Proposition 3.1 (Convergence Rate).**

When the functions  $f_n$  are convex, under assumption (A), and when  $\rho = L$ , we have

$$\mathbb{E}[f(\bar{\theta}_{n-1}) - f^*] \leq \frac{L\|\theta^* - \theta_0\|_2^2 + \frac{R^2}{L} \sum_{k=1}^n w_k^2}{2 \sum_{k=1}^n w_k} \quad \text{for all } n \geq 1, \quad (3)$$

where  $\bar{\theta}_{n-1}$  is defined in Algorithm 1,  $\theta^*$  is a minimizer of  $f$  on  $\Theta$ , and  $f^* \triangleq f(\theta^*)$ .

Such a rate is similar to the one of stochastic gradient descent with averaging, see [22] for example. Note that the constraint  $\rho = L$  here is compatible with the proximal gradient surrogate.

From Proposition 3.1, it is easy to obtain a  $O(1/\sqrt{n})$  bound for a finite horizon—that is, when the total number of iterations  $n$  is known in advance. When  $n$  is fixed, such a bound can indeed be obtained by plugging constant weights  $w_k = \gamma/\sqrt{n}$  for all  $k \leq n$  in Eq. (3). Note that the upper-bound  $O(1/\sqrt{n})$  cannot be improved in general without making further assumptions on the objective function [22]. The next corollary shows that in an infinite horizon setting and with decreasing weights, we lose a logarithmic factor compared to an optimal convergence rate [14, 22] of  $O(1/\sqrt{n})$ .

**Corollary 3.1 (Convergence Rate - Infinite Horizon - Decreasing Weights).**

Let us make the same assumptions as in Proposition 3.1 and choose the weights  $w_n = \gamma/\sqrt{n}$ . Then,

$$\mathbb{E}[f(\bar{\theta}_{n-1}) - f^*] \leq \frac{L\|\theta^* - \theta_0\|_2^2}{2\gamma\sqrt{n}} + \frac{R^2\gamma(1 + \log(n))}{2L\sqrt{n}}, \quad \forall n \geq 2.$$

Our analysis suggests to use weights of the form  $O(1/\sqrt{n})$ . In practice, we have found that choosing  $w_n = \sqrt{n_0 + 1}/\sqrt{n_0 + n}$  performs well, where  $n_0$  is tuned on a subsample of the training set.

### 3.2 Convergence Analysis - Strongly Convex Case

In this section, we introduce an additional assumption:

- (B) the functions  $f_n$  are  $\mu$ -strongly convex.

We show that our method achieves a rate  $O(1/n)$ , which is optimal up to a multiplicative constant for strongly convex functions (see [14, 22]).

**Proposition 3.2 (Convergence Rate).**

Under assumptions (A) and (B), with  $\rho = L + \mu$ . Define  $\beta \triangleq \frac{\mu}{\rho}$  and  $w_n \triangleq \frac{1+\beta}{1+\beta n}$ . Then,

$$\mathbb{E}[f(\hat{\theta}_{n-1}) - f^*] + \frac{\rho}{2} \mathbb{E}[\|\theta^* - \theta_n\|_2^2] \leq \max\left(\frac{2R^2}{\mu}, \rho\|\theta^* - \theta_0\|_2^2\right) \frac{1}{\beta n + 1} \quad \text{for all } n \geq 1,$$

where  $\hat{\theta}_n$  is defined in Algorithm 1, when choosing the averaging scheme called “option 3”.

The averaging scheme is slightly different than in the previous section and the weights decrease at a different speed. Again, this rate applies to the proximal gradient surrogates, which satisfy the constraint  $\rho = L + \mu$ . In the next section, we analyze our scheme in a non-convex setting.

### 3.3 Convergence Analysis - Non-Convex Case

Convergence results for non-convex problems are by nature weak, and difficult to obtain for stochastic optimization [4, 9]. In such a context, proving convergence to a global (or local) minimum is out of reach, and classical analyses study instead asymptotic stationary point conditions, which involve directional derivatives (see [2, 18]). Concretely, we introduce the following assumptions:

- (C)  $\Theta$  and the support  $\mathcal{X}$  of the data are compact;
- (D) The functions  $f_n$  are uniformly bounded by some constant  $M$ ;
- (E) The weights  $w_n$  are non-increasing,  $w_1 = 1$ ,  $\sum_{n \geq 1} w_n = +\infty$ , and  $\sum_{n \geq 1} w_n^2 \sqrt{n} < +\infty$ ;
- (F) The directional derivatives  $\nabla f_n(\theta, \theta' - \theta)$ , and  $\nabla f(\theta, \theta' - \theta)$  exist for all  $\theta$  and  $\theta'$  in  $\Theta$ .

Assumptions (C) and (D) combined with (A) are useful because they allow us to use some uniform convergence results from the theory of empirical processes [27]. In a nutshell, these assumptions ensure that the function class  $\{\mathbf{x} \mapsto \ell(\mathbf{x}, \theta) : \theta \in \Theta\}$  is “simple enough”, such that a uniform law of large numbers applies. The assumption (E) is more technical: it resembles classical conditions used for proving the convergence of stochastic gradient descent algorithms, usually stating that the weights  $w_n$  should be the summand of a diverging sum while the sum of  $w_n^2$  should be finite; the constraint  $\sum_{n \geq 1} w_n^2 \sqrt{n} < +\infty$  is slightly stronger. Finally, (F) is a mild assumption, which is useful to characterize the stationary points of the problem. A classical necessary first-order condition [2] for  $\theta$  to be a local minimum of  $f$  is indeed to have  $\nabla f(\theta, \theta' - \theta)$  non-negative for all  $\theta'$  in  $\Theta$ . We call such points  $\theta$  the stationary points of the function  $f$ . The next proposition is a generalization of a convergence result obtained in [19] in the context of sparse matrix factorization.

**Proposition 3.3 (Non-Convex Analysis - Almost Sure Convergence).**

Under assumptions (A), (C), (D), (E),  $(f(\theta_n))_{n \geq 0}$  converges with probability one. Under assumption (F), we also have that

$$\liminf_{n \rightarrow +\infty} \inf_{\theta \in \Theta} \frac{\nabla \bar{f}_n(\theta_n, \theta - \theta_n)}{\|\theta - \theta_n\|_2} \geq 0,$$

where the function  $\bar{f}_n$  is a weighted empirical risk recursively defined as  $\bar{f}_n = (1 - w_n)\bar{f}_{n-1} + w_n f_n$ . It can be shown that  $\bar{f}_n$  uniformly converges to  $f$ .

Even though  $\bar{f}_n$  converges uniformly to the expected cost  $f$ , Proposition 3.3 does not imply that the limit points of  $(\theta_n)_{n \geq 1}$  are stationary points of  $f$ . We obtain such a guarantee when the surrogates that are parameterized, an assumption always satisfied when Algorithm 1 is used in practice.

**Proposition 3.4 (Non-Convex Analysis - Parameterized Surrogates).**

Let us make the same assumptions as in Proposition 3.3, and let us assume that the functions  $\bar{g}_n$  are parameterized by some variables  $\kappa_n$  living in a compact set  $\mathcal{K}$  of  $\mathbb{R}^d$ . In other words,  $\bar{g}_n$  can be written as  $g_{\kappa_n}$ , with  $\kappa_n$  in  $\mathcal{K}$ . Suppose there exists a constant  $K > 0$  such that  $|g_{\kappa}(\theta) - g_{\kappa'}(\theta)| \leq K \|\kappa - \kappa'\|_2$  for all  $\theta$  in  $\Theta$  and  $\kappa, \kappa'$  in  $\mathcal{K}$ . Then, every limit point  $\theta_\infty$  of the sequence  $(\theta_n)_{n \geq 1}$  is a stationary point of  $f$ —that is, for all  $\theta$  in  $\Theta$ ,

$$\nabla f(\theta_\infty, \theta - \theta_\infty) \geq 0.$$

Finally, we show that our non-convex convergence analysis can be extended beyond first-order surrogate functions—that is, when  $g_n$  does not satisfy exactly Definition 2.1. This is possible when the objective has a particular partially separable structure, as shown in the next proposition. This extension was motivated by the non-convex sparse estimation formulation of Section 4, where such a structure appears.

**Proposition 3.5 (Non-Convex Analysis - Partially Separable Extension).**

Assume that the functions  $f_n$  split into  $f_n(\theta) = f_{0,n}(\theta) + \sum_{k=1}^K f_{k,n}(\gamma_k(\theta))$ , where the functions  $\gamma_k : \mathbb{R}^p \rightarrow \mathbb{R}$  are convex and  $R$ -Lipschitz, and the  $f_{k,n}$  are non-decreasing for  $k \geq 1$ . Consider  $g_{n,0}$  in  $\mathcal{S}_{L_0, \rho_1}(f_{0,n}, \theta_{n-1})$ , and some non-decreasing functions  $g_{k,n}$  in  $\mathcal{S}_{L_k, 0}(f_{k,n}, \gamma_k(\theta_{n-1}))$ . Instead of choosing  $g_n$  in  $\mathcal{S}_{L, \rho}(f_n, \theta_{n-1})$  in Alg 1, replace it by  $g_n \triangleq \theta \mapsto g_{0,n}(\theta) + g_{k,n}(\gamma_k(\theta))$ .

Then, Propositions 3.3 and 3.4 still hold.

## 4 Applications and Experimental Validation

In this section, we introduce different applications, and provide numerical experiments. A C++/Matlab implementation is available in the software package SPAMS [19].<sup>2</sup> All experiments were performed on a single core of a 2GHz Intel CPU with 64GB of RAM.

<sup>2</sup><http://spams-devel.gforge.inria.fr/>.

## 4.1 Stochastic Proximal Gradient Descent Algorithm

Our first application is a stochastic proximal gradient descent method, which we call SMM (Stochastic Majorization-Minimization), for solving problems of the form:

$$\min_{\theta \in \Theta} \mathbb{E}_{\mathbf{x}}[\ell(\mathbf{x}, \theta)] + \psi(\theta), \quad (4)$$

where  $\psi$  is a convex deterministic regularization function, and the functions  $\theta \mapsto \ell(\mathbf{x}, \theta)$  are differentiable and their gradients are  $L$ -Lipschitz continuous. We can thus use the proximal gradient surrogate presented in Section 2. Assume that a weight sequence  $(w_n)_{n \geq 1}$  is chosen such that  $w_1 = 1$ . By defining some other weights  $w_n^i$  recursively as  $w_n^i \triangleq (1 - w_n)w_n^{i-1}$  for  $i < n$  and  $w_n^n \triangleq w_n$ , our scheme yields the update rule:

$$\theta_n \leftarrow \arg \min_{\theta \in \Theta} \sum_{i=1}^n w_n^i \left[ \nabla f_i(\theta_{i-1})^\top \theta + \frac{L}{2} \|\theta - \theta_{i-1}\|_2^2 + \psi(\theta) \right]. \quad (\text{SMM})$$

Our algorithm is related to FOBOS [6], to SMIDAS [25] or the truncated gradient method [16] (when  $\psi$  is the  $\ell_1$ -norm). These three algorithms use indeed the following update rule:

$$\theta_n \leftarrow \arg \min_{\theta \in \Theta} \nabla f_n(\theta_{n-1})^\top \theta + \frac{1}{2\eta_n} \|\theta - \theta_{n-1}\|_2^2 + \psi(\theta), \quad (\text{FOBOS})$$

Another related scheme is the regularized dual averaging (RDA) of [30], which can be written as

$$\theta_n \leftarrow \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \nabla f_i(\theta_{i-1})^\top \theta + \frac{1}{2\eta_n} \|\theta\|_2^2 + \psi(\theta). \quad (\text{RDA})$$

Compared to these approaches, our scheme includes a weighted average of previously seen gradients, and a weighted average of the past iterates. Some links can also be drawn with approaches such as the ‘‘approximate follow the leader’’ algorithm of [10] and other works [12, 14].

We now evaluate the performance of our method for  $\ell_1$ -logistic regression. In summary, the datasets consist of pairs  $(y_i, \mathbf{x}_i)_{i=1}^N$ , where the  $y_i$ ’s are in  $\{-1, +1\}$ , and the  $\mathbf{x}_i$ ’s are in  $\mathbb{R}^p$  with unit  $\ell_2$ -norm. The function  $\psi$  in (4) is the  $\ell_1$ -norm:  $\psi(\theta) \triangleq \lambda \|\theta\|_1$ , and  $\lambda$  is a regularization parameter; the functions  $f_i$  are logistic losses:  $f_i(\theta) \triangleq \log(1 + e^{-y_i \mathbf{x}_i^\top \theta})$ . One part of each dataset is devoted to training, and another part to testing. We used weights of the form  $w_n \triangleq \sqrt{(n_0 + 1)/(n + n_0)}$ , where  $n_0$  is automatically adjusted at the beginning of each experiment by performing one pass on 5% of the training data. We implemented SMM in C++ and exploited the sparseness of the datasets, such that each update has a computational complexity of the order  $O(s)$ , where  $s$  is the number of non zeros in  $\nabla f_n(\theta_{n-1})$ ; such an implementation is non trivial but proved to be very efficient.

We consider three datasets described in the table below. `rcv1` and `webspam` are obtained from the 2008 Pascal large-scale learning challenge.<sup>3</sup> `kdd2010` is available from the LIBSVM website.<sup>4</sup>

name	$N_{\text{tr}}$ (train)	$N_{\text{te}}$ (test)	$p$	density (%)	size (GB)
<code>rcv1</code>	781 265	23 149	47 152	0.161	0.95
<code>webspam</code>	250 000	100 000	16 091 143	0.023	14.95
<code>kdd2010</code>	10 000 000	9 264 097	28 875 157	$10^{-4}$	4.8

We compare our implementation with state-of-the-art publicly available solvers: the batch algorithm FISTA of [1] implemented in the C++ SPAMS toolbox and LIBLINEAR v1.93 [7]. LIBLINEAR is based on a working-set algorithm, and, to the best of our knowledge, is one of the most efficient available solver for  $\ell_1$ -logistic regression with sparse datasets. Because  $p$  is large, the incremental majorization-minimization method of [18] could not run for memory reasons. We run every method on 1, 2, 3, 4, 5, 10 and 25 epochs (passes over the training set), for three regularization regimes, respectively yielding a solution with approximately 100, 1 000 and 10 000 non-zero coefficients. We report results for the medium regularization in Figure 1 and provide the rest as supplemental material. FISTA is not represented in this figure since it required more than 25 epochs to achieve reasonable values. Our conclusion is that *SMM often provides a reasonable solution after one epoch, and outperforms LIBLINEAR in the low-precision regime. For high-precision regimes, LIBLINEAR should be preferred.* Such a conclusion is often obtained when comparing batch and stochastic algorithms [4], but matching the performance of LIBLINEAR is very challenging.

<sup>3</sup><http://largescale.ml.tu-berlin.de>.

<sup>4</sup><http://www.csie.ntu.edu.tw/~cjlin/libsvm/>.

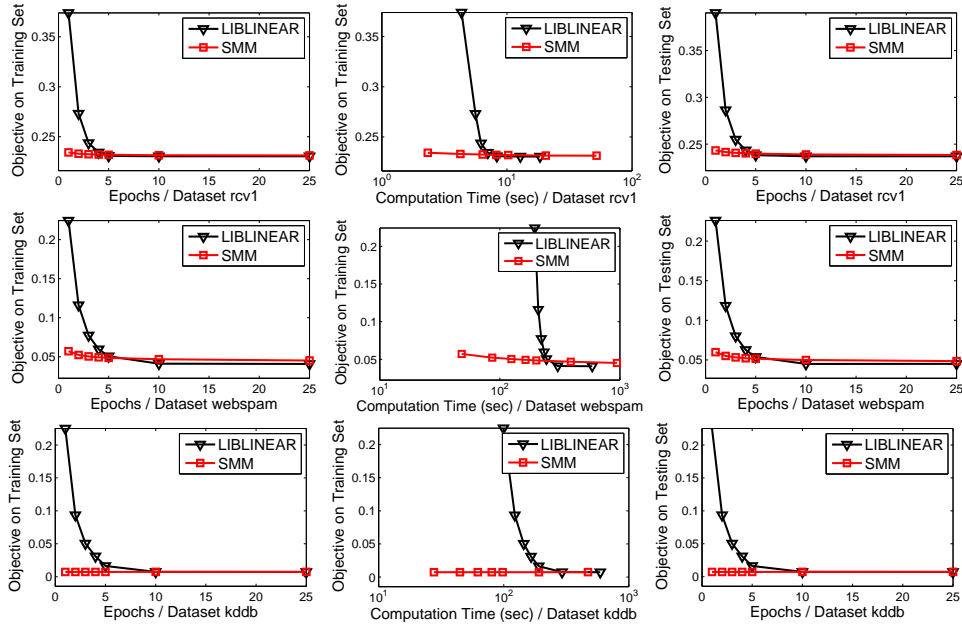


Figure 1: Comparison between LIBLINEAR and SMM for the medium regularization regime.

## 4.2 Online DC Programming for Non-Convex Sparse Estimation

We now consider the same experimental setting as in the previous section, but with a non-convex regularizer  $\psi : \theta \mapsto \lambda \sum_{j=1}^p \log(|\theta[j]| + \varepsilon)$ , where  $\theta[j]$  is the  $j$ -th entry in  $\theta$ . A classical way for minimizing the regularized empirical cost  $\frac{1}{N} \sum_{i=1}^N f_i(\theta) + \psi(\theta)$  is to resort to DC programming. It consists of solving a sequence of reweighted- $\ell_1$  problems [8]. A current estimate  $\theta_{n-1}$  is updated as a solution of  $\min_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N f_i(\theta) + \lambda \sum_{j=1}^p \eta_j |\theta[j]|$ , where  $\eta_j \triangleq 1/(|\theta_{n-1}[j]| + \varepsilon)$ .

In contrast to this “batch” methodology, we can use our framework to address the problem online. At iteration  $n$  of Algorithm 1, we define the function  $g_n$  according to Proposition 3.5:

$$g_n : \theta \mapsto f_n(\theta_{n-1}) + \nabla f_n(\theta_{n-1})^\top (\theta - \theta_{n-1}) + \frac{L}{2} \|\theta - \theta_{n-1}\|_2^2 + \lambda \sum_{j=1}^p \frac{|\theta[j]|}{|\theta_{n-1}[j]| + \varepsilon},$$

We compare our online DC programming algorithm against the batch one, and report the results in Figure 2, with  $\varepsilon$  set to 0.01. We conclude that *the batch reweighted- $\ell_1$  algorithm always converges after 2 or 3 weight updates, but suffers from local minima issues. The stochastic algorithm exhibits a slower convergence, but provides significantly better solutions.* Whether or not there are good theoretical reasons for this fact remains to be investigated. Note that it would have been more rigorous to choose a bounded set  $\Theta$ , which is required by Proposition 3.5. In practice, it turns not to be necessary for our method to work well; the iterates  $\theta_n$  have indeed remained in a bounded set.

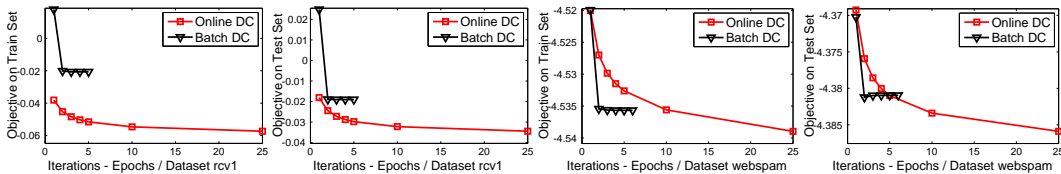


Figure 2: Comparison between batch and online DC programming, with medium regularization for the datasets rcv1 and webspam. Additional plots are provided in the supplemental material. Note that each iteration in the batch setting can perform several epochs (passes over training data).

## 4.3 Online Structured Sparse Coding

In this section, we show that we can bring new functionalities to existing matrix factorization techniques [13, 19]. We are given a large collection of signals  $(\mathbf{x}_i)_{i=1}^N$  in  $\mathbb{R}^m$ , and we want to find a



dictionary  $\mathbf{D}$  in  $\mathbb{R}^{m \times K}$  that can represent these signals in a sparse way. The quality of  $\mathbf{D}$  is measured through the loss  $\ell(\mathbf{x}, \mathbf{D}) \triangleq \min_{\alpha \in \mathbb{R}^K} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\alpha\|_2^2 + \lambda_1 \|\alpha\|_1 + \frac{\lambda_2}{2} \|\alpha\|_2^2$ , where the  $\ell_1$ -norm can be replaced by any convex regularizer, and the squared loss by any convex smooth loss.

Then, we are interested in minimizing the following expected cost:

$$\min_{\mathbf{D} \in \mathbb{R}^{m \times K}} \mathbb{E}_{\mathbf{x}} [\ell(\mathbf{x}, \mathbf{D})] + \varphi(\mathbf{D}),$$

where  $\varphi$  is a regularizer for  $\mathbf{D}$ . In the online learning approach of [19], the only way to regularize  $\mathbf{D}$  is to use a constraint set, on which we need to be able to project efficiently; this is unfortunately not always possible. In the matrix factorization framework of [13], it is argued that some applications can benefit from a structured penalty  $\varphi$ , but the approach of [13] is not easily amenable to stochastic optimization. Our approach makes it possible by using the proximal gradient surrogate

$$g_n : \mathbf{D} \mapsto \ell(\mathbf{x}_n, \mathbf{D}_{n-1}) + \text{Tr} (\nabla_{\mathbf{D}} \ell(\mathbf{x}_n, \mathbf{D}_{n-1})^\top (\mathbf{D} - \mathbf{D}_{n-1})) + \frac{\ell}{2} \|\mathbf{D} - \mathbf{D}_{n-1}\|_F^2 + \varphi(\mathbf{D}). \quad (5)$$

It is indeed possible to show that  $\mathbf{D} \mapsto \ell(\mathbf{x}_n, \mathbf{D})$  is differentiable, and its gradient is Lipschitz continuous with a constant  $L$  that can be explicitly computed [18, 19].

We now design a proof-of-concept experiment. We consider a set of  $N = 400\,000$  whitened natural image patches  $\mathbf{x}_n$  of size  $m = 20 \times 20$  pixels. We visualize some elements from a dictionary  $\mathbf{D}$  trained by SPAMS [19] on the left of Figure 3; the dictionary elements are almost sparse, but have some residual noise among the small coefficients. Following [13], we propose to use a regularization function  $\varphi$  encouraging neighbor pixels to be set to zero together, thus leading to a sparse structured dictionary. We consider the collection  $\mathcal{G}$  of all groups of variables corresponding to squares of 4 neighbor pixels in  $\{1, \dots, m\}$ . Then, we define  $\varphi(\mathbf{D}) \triangleq \gamma_1 \sum_{j=1}^K \sum_{g \in \mathcal{G}} \max_{k \in g} |\mathbf{d}_j[k]| + \gamma_2 \|\mathbf{D}\|_F^2$ , where  $\mathbf{d}_j$  is the  $j$ -th column of  $\mathbf{D}$ . The penalty  $\varphi$  is a structured sparsity-inducing penalty that encourages groups of variables  $g$  to be set to zero together [13]. Its proximal operator can be computed efficiently [20], and it is thus easy to use the surrogates (5). We set  $\lambda_1 = 0.15$  and  $\lambda_2 = 0.01$ ; after trying a few values for  $\gamma_1$  and  $\gamma_2$  at a reasonable computational cost, we obtain dictionaries with the desired regularization effect, as shown in Figure 3. Learning one dictionary of size  $K = 256$  took a few minutes when performing one pass on the training data with mini-batches of size 100. This experiment demonstrates that our approach is more flexible and general than [13] and [19]. Note that it is possible to show that when  $\gamma_2$  is large enough, the iterates  $\mathbf{D}_n$  necessarily remain in a bounded set, and thus our convergence analysis presented in Section 3.3 applies to this experiment.

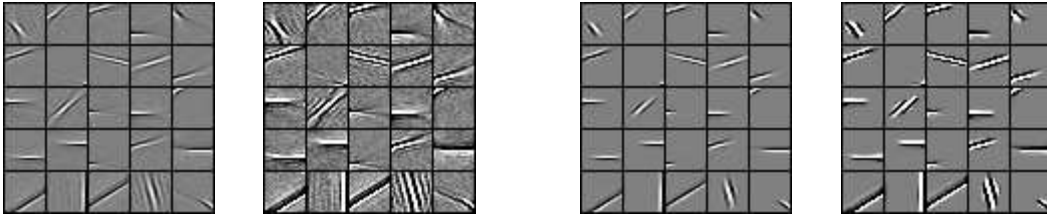


Figure 3: Left: Two visualizations of 25 elements from a larger dictionary obtained by the toolbox SPAMS [19]; the second view amplifies the small coefficients. Right: the corresponding views of the dictionary elements obtained by our approach after initialization with the dictionary on the left.

## 5 Conclusion

In this paper, we have introduced a stochastic majorization-minimization algorithm that gracefully scales to millions of training samples. We have shown that it has strong theoretical properties and some practical value in the context of machine learning. We have derived from our framework several new algorithms, which have shown to match or outperform the state of the art for solving large-scale convex problems, and to open up new possibilities for non-convex ones. In the future, we would like to study surrogate functions that can exploit the curvature of the objective function, which we believe is a crucial issue to deal with badly conditioned datasets.

## Acknowledgments

This work was supported by the Gargantua project (program Mastodons - CNRS).

## References

- [1] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.*, 2(1):183–202, 2009.
- [2] J.M. Borwein and A.S. Lewis. *Convex analysis and nonlinear optimization*. Springer, 2006.
- [3] L. Bottou. Online algorithms and stochastic approximations. In David Saad, editor, *Online Learning and Neural Networks*. 1998.
- [4] L. Bottou and O. Bousquet. The trade-offs of large scale learning. In *Adv. NIPS*, 2008.
- [5] O. Cappé and E. Moulines. On-line expectation–maximization algorithm for latent data models. *J. Roy. Stat. Soc. B*, 71(3):593–613, 2009.
- [6] J. Duchi and Y. Singer. Efficient online and batch learning using forward backward splitting. *J. Mach. Learn. Res.*, 10:2899–2934, 2009.
- [7] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. LIBLINEAR: A library for large linear classification. *J. Mach. Learn. Res.*, 9:1871–1874, 2008.
- [8] G. Gasso, A. Rakotomamonjy, and S. Canu. Recovering sparse signals with non-convex penalties and DC programming. *IEEE T. Signal Process.*, 57(12):4686–4698, 2009.
- [9] S. Ghadimi and G. Lan. Stochastic first- and zeroth-order methods for nonconvex stochastic programming. Technical report, 2013.
- [10] E. Hazan, A. Agarwal, and S. Kale. Logarithmic regret algorithms for online convex optimization. *Mach. Learn.*, 69(2-3):169–192, 2007.
- [11] E. Hazan and S. Kale. Beyond the regret minimization barrier: an optimal algorithm for stochastic strongly-convex optimization. In *Proc. COLT*, 2011.
- [12] C. Hu, J. Kwok, and W. Pan. Accelerated gradient methods for stochastic optimization and online learning. In *Adv. NIPS*, 2009.
- [13] R. Jenatton, G. Obozinski, and F. Bach. Structured sparse principal component analysis. In *Proc. AIS-TATS*, 2010.
- [14] G. Lan. An optimal method for stochastic composite optimization. *Math. Program.*, 133:365–397, 2012.
- [15] K. Lange, D.R. Hunter, and I. Yang. Optimization transfer using surrogate objective functions. *J. Comput. Graph. Stat.*, 9(1):1–20, 2000.
- [16] J. Langford, L. Li, and T. Zhang. Sparse online learning via truncated gradient. *J. Mach. Learn. Res.*, 10:777–801, 2009.
- [17] N. Le Roux, M. Schmidt, and F. Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In *Adv. NIPS*, 2012.
- [18] J. Mairal. Optimization with first-order surrogate functions. In *Proc. ICML*, 2013. arXiv:1305.3120.
- [19] J. Mairal, F. Bach, J. Ponce, and G. Sapiro. Online learning for matrix factorization and sparse coding. *J. Mach. Learn. Res.*, 11:19–60, 2010.
- [20] J. Mairal, R. Jenatton, G. Obozinski, and F. Bach. Network flow algorithms for structured sparsity. In *Adv. NIPS*, 2010.
- [21] R.M. Neal and G.E. Hinton. A view of the EM algorithm that justifies incremental, sparse, and other variants. *Learning in graphical models*, 89, 1998.
- [22] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM J. Optimiz.*, 19(4):1574–1609, 2009.
- [23] Y. Nesterov. Gradient methods for minimizing composite objective functions. Technical report, CORE Discussion Paper, 2007.
- [24] S. Shalev-Schwartz and T. Zhang. Proximal stochastic dual coordinate ascent. *arXiv 1211.2717v1*, 2012.
- [25] S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Stochastic convex optimization. In *Proc. COLT*, 2009.
- [26] S. Shalev-Shwartz and A. Tewari. Stochastic methods for  $\ell_1$  regularized loss minimization. In *Proc. ICML*, 2009.
- [27] A. W. Van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.
- [28] M.J. Wainwright and M.I. Jordan. Graphical models, exponential families, and variational inference. *Found. Trends Mach. Learn.*, 1(1-2):1–305, 2008.
- [29] S. Wright, R. Nowak, and M. Figueiredo. Sparse reconstruction by separable approximation. *IEEE T. Signal Process.*, 57(7):2479–2493, 2009.
- [30] L. Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *J. Mach. Learn. Res.*, 11:2543–2596, 2010.

## A Mathematical Background and Useful Results

In this paper, we use subdifferential calculus for convex functions. The definition of subgradients and directional derivatives can be found in classical textbooks, see, e.g., [2], [37]. We denote by  $\partial f(\theta)$  the subdifferential of a convex function  $f$  at a point  $\theta$ . Other definitions can be found in the appendix of [18], which uses a similar notation as ours.

In this section, we present several classical optimization and probabilistic tools, which we use in our paper. The first lemma is a classical quadratic upper-bound for differentiable functions with a Lipschitz gradient. It can be found for instance in Lemma 1.2.3 of [35], or in the appendix of [18].

**Lemma A.1 (Convex Surrogate for Functions with Lipschitz Gradient).**

Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be differentiable and  $\nabla f$  be  $L$ -Lipschitz continuous. Then, for all  $\theta, \theta'$  in  $\mathbb{R}^p$ ,

$$|f(\theta') - f(\theta) - \nabla f(\theta)^\top (\theta' - \theta)| \leq \frac{L}{2} \|\theta - \theta'\|_2^2. \quad (6)$$

The next lemma is a simple relation, which will allow us to identify the subdifferential of a convex function with the one of its surrogate at a particular point.

**Lemma A.2 (Surrogate Functions and Subdifferential).**

Assume that  $f, g : \mathbb{R}^p \rightarrow \mathbb{R}$  are convex, and that  $h \triangleq g - f$  is differentiable at  $\theta$  in  $\mathbb{R}^p$  with  $\nabla h(\theta) = 0$ . Then,  $\partial f(\theta) = \partial g(\theta)$ .

*Proof.* It is easy to show that  $g$  and  $f$  have the same directional derivatives at  $\theta$  since  $h$  is differentiable and  $\nabla h(\theta) = 0$ . This is sufficient to conclude that  $\partial g(\theta) = \partial f(\theta)$  by using Proposition 3.1.6 of [2], a simple lemma relating directional derivatives and subgradients.  $\square$

The following lemma is a lower bound for strongly convex functions. It can be found for instance in [36].

**Lemma A.3 (Lower Bound for Strongly Convex Functions).**

Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a  $\mu$ -strongly convex function. Let  $z$  be in  $\partial f(\kappa)$  for some  $\kappa$  in  $\mathbb{R}^p$ . Then, the following inequality holds for all  $\theta$  in  $\mathbb{R}^p$ :

$$f(\theta) \geq f(\kappa) + z^\top (\theta - \kappa) + \frac{\mu}{2} \|\theta - \kappa\|_2^2.$$

*Proof.* The function  $l : \theta \mapsto f(\theta) - \frac{\mu}{2} \|\theta - \kappa\|_2^2$  is convex by definition of strong convexity, and  $l - f$  is differentiable with  $\nabla(l - f)(\kappa) = 0$ . We apply Lemma A.2, which tells us that  $z$  is in  $\partial l(\kappa)$ . This is sufficient to conclude, by noticing that a convex function is always above its tangents.  $\square$

The next lemma is also classical (see the appendix of [18]).

**Lemma A.4 (Second-Order Growth Property).**

Let  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  be a  $\mu$ -strongly convex function and  $\Theta \subseteq \mathbb{R}^p$  be a convex set. Let  $\theta^*$  be the minimizer of  $f$  on  $\Theta$ . Then, the following condition holds for all  $\theta$  in  $\Theta$ :

$$f(\theta) \geq f(\theta^*) + \frac{\mu}{2} \|\theta - \theta^*\|_2^2.$$

We now introduce a sequence of probabilistic tools, which we use in our convergence analysis for non-convex functions. The first one is a classical theorem on quasi-martingales, which was used in [3] for proving the convergence of the stochastic gradient descent algorithm.

**Theorem A.1 (Convergence of Quasi-Martingales).**

This presentation follows [3] and Proposition 9.5 and Theorem 9.4 of [34]. The original theorem is due to [33]. Let  $(\mathcal{F}_n)_{n \geq 0}$  be an increasing family of  $\sigma$ -fields. Let  $(X_n)_{n \geq 0}$  be a real stochastic process such that every random variable  $X_n$  is bounded below by a constant independent of  $n$ , and  $\mathcal{F}_n$ -measurable. Let

$$\delta_n \triangleq \begin{cases} 1 & \text{if } \mathbb{E}[X_{n+1} - X_n | \mathcal{F}_n] > 0, \\ 0 & \text{otherwise.} \end{cases}$$

If the series  $\sum_{n=0}^{\infty} \mathbb{E}[\delta_n(X_{n+1} - X_n)]$  converges, then  $(X_n)_{n \geq 0}$  is a quasi-martingale and converges almost surely to an integrable random variable  $X_{\infty}$ . Moreover,

$$\sum_{n=0}^{\infty} \mathbb{E}[|\mathbb{E}[X_{n+1} - X_n | \mathcal{F}_n]|] < \infty.$$

The next lemma is simple, but useful to prove the convergence of deterministic algorithms.

**Lemma A.5. Deterministic Lemma on Non-negative Converging Series.**

Let  $(a_n)_{n \geq 1}$ ,  $(b_n)_{n \geq 1}$  be two non-negative real sequences such that the series  $\sum_{n=1}^{\infty} a_n$  diverges, the series  $\sum_{n=1}^{\infty} a_n b_n$  converges, and there exists  $K > 0$  such that  $|b_{n+1} - b_n| \leq K a_n$ . Then, the sequence  $(b_n)_{n \geq 1}$  converges to 0.

*Proof.* The proof is inspired by the one of Proposition 1.2.4 of [31]. Since the series  $\sum_{n \geq 1} a_n$  diverges, we necessarily have  $\liminf_{n \rightarrow +\infty} b_n = 0$ . Otherwise, it would be easy to contradict the assumption  $\sum_{n \geq 1} a_n b_n < +\infty$ .

Let us now proceed by contradiction and assume that  $\limsup_{n \rightarrow +\infty} b_n = \lambda > 0$ . We can then build two sequences of indices  $(m_j)_{j \geq 1}$  and  $(n_j)_{j \geq 1}$  such that

- $m_j < n_j < m_{j+1}$ ,
- $\frac{\lambda}{3} < b_k$ , for  $m_j \leq k < n_j$ ,
- $b_k \leq \frac{\lambda}{3}$ , for  $n_j \leq k < m_{j+1}$ .

Let  $\varepsilon = \frac{\lambda^2}{9K}$  and  $\tilde{j}$  be large enough such that

$$\sum_{n=m_j}^{\infty} a_n b_n < \varepsilon.$$

Then, we have for all  $j \geq \tilde{j}$  and all  $m$  with  $m_j \leq m \leq n_j - 1$ ,

$$\begin{aligned} |b_{n_j} - b_m| &\leq \sum_{k=m}^{n_j-1} |b_{k+1} - b_k| \leq \frac{3K}{\lambda} \sum_{k=m}^{n_j-1} a_k \frac{\lambda}{3} \leq \frac{3K}{\lambda} \sum_{k=m}^{n_j-1} a_k b_k \leq \frac{3K}{\lambda} \sum_{k=m}^{+\infty} a_k b_k \\ &\leq \frac{3K\varepsilon}{\lambda} \leq \frac{\lambda}{3}. \end{aligned}$$

Therefore, by using the triangle inequality,

$$b_m \leq b_{n_j} + \frac{\lambda}{3} \leq \frac{2\lambda}{3}.$$

and finally, for all  $m \geq \tilde{j}$ ,

$$b_m \leq \frac{2\lambda}{3},$$

which contradicts  $\limsup_{n \rightarrow +\infty} b_n = \lambda > 0$ . Therefore,  $b_n \xrightarrow{n \rightarrow +\infty} 0$ .

□

We now provide a stochastic version of Lemma A.6.

**Lemma A.6. Stochastic Lemma on Non-negative Converging Series.**

Let  $(X_n)_{n \geq 1}$  be a sequence of non-negative measurable random variables on a probability space. Let also  $a_n, b_n$  be two non-negative sequences such that  $\sum_{n \geq 1} a_n = +\infty$  and  $\sum_{n \geq 1} a_n b_n < +\infty$ . Assume that there exists a constant  $C$  such that for all  $n \geq 1$ ,  $\mathbb{E}[X_n] \leq b_n$  and  $|X_{n+1} - X_n| \leq C a_n$  almost surely. Then  $X_n$  almost surely converges to zero.

*Proof.* The following series is convergent

$$\mathbb{E} \left[ \sum_{n \geq 1} a_n X_n \right] = \sum_{n \geq 1} \mathbb{E} [a_n X_n] \leq \sum_{n \geq 1} a_n b_n < +\infty,$$

where we use the fact that the random variables are non-negative to interchange the sum and the expectation. We thus have that  $\sum_{n \geq 1} a_n X_n$  converges with probability one. Then, let us call  $a'_n = a_n$  and  $b'_n = X_n$ ; the conditions of Lemma A.5 are satisfied for  $a'_n$  and  $b'_n$  with probability one, and  $X_n$  almost surely converges to zero.  $\square$

## B Auxiliary Lemmas

In this section, we present auxiliary lemmas for our convex and non-convex analyses. We start by presenting a lemma which is useful for both of them, and which is in fact a core component for all results presented in [18]. The proof of this lemma is simple and available in [18].

### Lemma B.1 (Basic Properties of First-Order Surrogate Functions).

Let  $g$  be in  $\mathcal{S}_{L,\rho}(f, \kappa)$  for some  $\kappa$  in  $\Theta$ . Define the approximation error function  $h \triangleq g - f$  and let  $\theta'$  be the minimizer of  $g$  over  $\Theta$ . Then, for all  $\theta$  in  $\Theta$ ,

- $\nabla h(\kappa) = 0$ ;
- $|h(\theta)| \leq \frac{L}{2} \|\theta - \kappa\|_2^2$ ;
- $f(\theta') \leq g(\theta') \leq f(\theta) + \frac{L}{2} \|\theta - \kappa\|_2^2 - \frac{\rho}{2} \|\theta - \theta'\|_2^2$ .

### B.1 Convex Analysis

We introduce, for all  $n \geq 0$ , the quantity  $\xi_n \triangleq \frac{1}{2} \mathbb{E}[\|\theta^* - \theta_n\|_2^2]$ , where  $\theta^*$  is a minimizer of  $f$  on  $\Theta$ . Our analysis also involves several quantities that are defined recursively for all  $n \geq 1$ :

$$\begin{cases} A_n \triangleq (1 - w_n)A_{n-1} + w_n \xi_{n-1} \\ B_n \triangleq (1 - w_n)B_{n-1} + w_n \mathbb{E}[f(\theta_{n-1})] \\ C_n \triangleq (1 - w_n)C_{n-1} + \frac{(Rw_n)^2}{2\rho} \\ \bar{g}_n \triangleq (1 - w_n)\bar{g}_{n-1} + w_n g_n \\ \bar{f}_n \triangleq (1 - w_n)\bar{f}_{n-1} + w_n f_n \end{cases}, \quad (7)$$

where  $A_0 \triangleq \frac{1}{L}(\rho\xi_0 - f^*)$ ,  $B_0 \triangleq 0$ ,  $C_0 \triangleq 0$ ,  $\bar{g}_0 = \bar{f}_0 \triangleq \theta \mapsto \frac{\rho}{2} \|\theta - \theta_0\|_2^2$ . Note that  $\bar{g}_0$  is  $\rho$ -strongly convex, and is minimized by  $\theta_0$ . The choice for  $A_0, B_0, C_0$  is driven by technical reasons, which appear in the proof of Lemma B.4, a stochastic version of Lemma B.1.

Note that we also assume here that all the expectations above are well defined and finite-valued. In other words, we do not deal with measurability or integrability issues for simplicity, as often done in the literature [22].

### Lemma B.2 (Auxiliary Lemma for Convex Analysis).

When the functions  $f_n$  are convex, and the surrogates  $g_n$  are in  $\mathcal{S}_{L,\rho}(f_n, \theta_{n-1})$ , we have under assumption (A) that for all  $n \geq 1$ ,

$$\bar{g}_n(\theta_{n-1}) \leq \bar{g}_n(\theta_n) + \frac{(Rw_n)^2}{2\rho}.$$

*Proof.* First, we remark that the subdifferentials of  $g_n$  and  $f_n$  at  $\theta_{n-1}$  coincide by applying Lemma A.2. Then, we choose  $z_n$  in  $\partial g_n(\theta_{n-1}) = \partial f_n(\theta_{n-1})$ , which is bounded by  $R$  accord-

ing to assumption **(A)**, and we have

$$\begin{aligned}
\bar{g}_n(\theta_n) &= (1 - w_n)\bar{g}_{n-1}(\theta_n) + w_n g_n(\theta_n) \\
&\geq (1 - w_n) \left( \bar{g}_{n-1}(\theta_{n-1}) + \frac{\rho}{2} \|\theta_n - \theta_{n-1}\|_2^2 \right) + w_n \left( g_n(\theta_{n-1}) + z_n^\top (\theta_n - \theta_{n-1}) + \frac{\rho}{2} \|\theta_n - \theta_{n-1}\|_2^2 \right) \\
&= \bar{g}_n(\theta_{n-1}) + w_n z_n^\top (\theta_n - \theta_{n-1}) + \frac{\rho}{2} \|\theta_n - \theta_{n-1}\|_2^2 \\
&\geq \bar{g}_n(\theta_{n-1}) - R w_n \|\theta_n - \theta_{n-1}\|_2 + \frac{\rho}{2} \|\theta_n - \theta_{n-1}\|_2^2 \\
&\geq \bar{g}_n(\theta_{n-1}) - \frac{(R w_n)^2}{2\rho}.
\end{aligned}$$

The first inequality uses Lemma A.4 and Lemma A.3 since  $g_n$  is  $\rho$ -strongly convex by definition (and by induction  $\bar{g}_n$  is  $\rho$ -strongly convex as well); the second inequality uses Cauchy-Schwarz's inequality and the fact that the subgradients of the functions  $f_n$  are bounded by  $R$ .  $\square$

**Lemma B.3 (Another Auxiliary Lemma for Convex Analysis).**

When the functions  $f_n$  are convex, and the surrogates  $g_n$  are in  $\mathcal{S}_{L,\rho}(f_n, \theta_{n-1})$ , we have under assumption **(A)** that for all  $n \geq 0$ ,

$$B_n \leq \mathbb{E}[\bar{g}_n(\theta_n)] + C_n, \quad (8)$$

*Proof.* We proceed by induction, and start by showing that Eq. (8) is true for  $n = 0$ .

$$B_0 = 0 = \mathbb{E}[\bar{g}_0(\theta_0)] = \mathbb{E}[\bar{g}_0(\theta_0)] + C_0.$$

Let us now assume that it is true for  $n - 1$ , and show that it is true for  $n$ .

$$\begin{aligned}
B_n &= (1 - w_n)B_{n-1} + w_n \mathbb{E}[f(\theta_{n-1})] \\
&\leq (1 - w_n)(\mathbb{E}[\bar{g}_{n-1}(\theta_{n-1})] + C_{n-1}) + w_n \mathbb{E}[f(\theta_{n-1})] \\
&= (1 - w_n)(\mathbb{E}[\bar{g}_{n-1}(\theta_{n-1})] + C_{n-1}) + w_n \mathbb{E}[f_n(\theta_{n-1})] \\
&= (1 - w_n)(\mathbb{E}[\bar{g}_{n-1}(\theta_{n-1})] + C_{n-1}) + w_n \mathbb{E}[g_n(\theta_{n-1})] \\
&= \mathbb{E}[\bar{g}_n(\theta_{n-1})] + (1 - w_n)C_{n-1} \\
&\leq \mathbb{E}[\bar{g}_n(\theta_n)] + \frac{(R w_n)^2}{2\rho} + (1 - w_n)C_{n-1} \\
&= \mathbb{E}[\bar{g}_n(\theta_n)] + C_n.
\end{aligned}$$

The first inequality uses the induction hypothesis; the last inequality uses Lemma B.2 and the definition of  $C_n$ . We also used the fact that  $\mathbb{E}[f_n(\theta_{n-1})] = \mathbb{E}[\mathbb{E}[f_n(\theta_{n-1})|\mathcal{F}_{n-1}]] = \mathbb{E}[\mathbb{E}[f(\theta_{n-1})|\mathcal{F}_{n-1}]] = \mathbb{E}[f(\theta_{n-1})]$ , where  $\mathcal{F}_{n-1}$  corresponds to the filtration induced by the past information before time  $n$ , such that  $\theta_{n-1}$  is deterministic given  $\mathcal{F}_{n-1}$ .  $\square$

The next lemma is important; it is the stochastic version of Lemma B.1 for first-order surrogates.

**Lemma B.4 (Basic Properties of Stochastic First-Order Surrogates).**

When the functions  $f_n$  are convex and the functions  $g_n$  are in  $\mathcal{S}_{L,\rho}(f_n, \theta_{n-1})$ , we have under assumption **(A)** that for all  $n \geq 0$ ,

$$B_n \leq f^* + L A_n - \rho \xi_n + C_n,$$

*Proof.* According to Lemma B.3, it is sufficient to show that  $\mathbb{E}[\bar{g}_n(\theta_n)] \leq f^* + L A_n - \rho \xi_n$  for all  $n \geq 0$ . Since  $\bar{g}_n$  is  $\rho$ -strongly convex and  $\theta_n$  is the minimizer of  $\bar{g}_n$  over  $\Theta$ , we have  $\mathbb{E}[\bar{g}_n(\theta_n)] \leq \mathbb{E}[\bar{g}_n(\theta^*)] - \rho \xi_n$ , by using Lemma A.4. Thus, it is in fact sufficient to show that  $\mathbb{E}[\bar{g}_n(\theta^*)] \leq f^* + L A_n$ . For  $n = 0$ , this inequality holds since  $\mathbb{E}[\bar{g}_0(\theta^*)] = \rho \xi_0 = f^* + L A_0$ . We can then proceed again by induction: assume that  $\mathbb{E}[\bar{g}_{n-1}(\theta^*)] \leq f^* + L A_{n-1}$ . Then,

$$\begin{aligned}
\mathbb{E}[\bar{g}_n(\theta^*)] &= (1 - w_n)\mathbb{E}[\bar{g}_{n-1}(\theta^*)] + w_n \mathbb{E}[g_n(\theta^*)] \\
&\leq (1 - w_n)(f^* + L A_{n-1}) + w_n (\mathbb{E}[f_n(\theta^*)] + L \xi_{n-1}) \\
&= (1 - w_n)(f^* + L A_{n-1}) + w_n (f^* + L \xi_{n-1}) \\
&= f^* + L A_n,
\end{aligned}$$

where we have used Lemma B.1 to upper-bound the difference  $\mathbb{E}[g_n(\theta^*)] - \mathbb{E}[f_n(\theta^*)]$  by  $\xi_{n-1}$ .  $\square$

For strongly-convex functions, we also have the following simple but useful relation between  $A_n$  and  $B_n$ .

**Lemma B.5 (Relation between  $A_n$  and  $B_n$ ).**

Under assumption **(B)**, if  $w_1 = 1$ , we have for all  $n \geq 1$ ,

$$f^* + \mu A_n \leq B_n.$$

*Proof.* This relation is true for  $n = 1$  since we have  $f^* + \mu A_1 = f^* + \mu \xi_0 \leq f(\theta_0) = B_1$  by applying Lemma A.4, since  $f$  is  $\mu$ -strongly convex according to assumption **(B)**. The rest follows by induction.  $\square$

## B.2 Non-convex Analysis

When the functions  $f_n$  are not convex, the convergence analysis becomes more involved. One key tool we use is a uniform convergence result when the function class  $\{\mathbf{x} \mapsto \ell(\theta, \mathbf{x}) : \theta \in \Theta\}$  is “simple enough” in terms of entropy. Under the assumptions made in our paper, it is indeed possible to use some results from empirical processes [27], which provides us the following lemma.

**Lemma B.6 (Uniform Convergence).**

Under assumptions **(A)**, **(C)**, and **(D)**, we have the following uniform law of large numbers:

$$\mathbb{E} \left[ \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n f_i(\theta) - f(\theta) \right| \right] \leq \frac{C}{\sqrt{n}}, \quad (9)$$

where  $C$  is a constant, and  $\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n f_i(\theta) - f(\theta) \right|$  converges almost surely to zero.

*Proof.* We simply refer to Lemma 19.36 and Example 19.7 of [27], where assumptions **(C)** and **(D)** ensure uniform boundness and squared integrability conditions. Note that we assume that the quantities  $\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n f_i(\theta) - f(\theta) \right|$  are measurable. This assumption does not incur a loss of generality, since measurability issues for empirical processes can be dealt with rigorously [27].  $\square$

The next lemma shows that uniform convergence applies to the weighted empirical risk  $\bar{f}_n$ , defined in Eq. (7), but with a different rate.

**Lemma B.7 (Uniform Convergence for  $\bar{f}_n$ ).**

Under assumptions **(A)**, **(C)**, **(D)**, and **(E)**, we have for all  $n \geq 1$ ,

$$\mathbb{E} \left[ \sup_{\theta \in \Theta} |\bar{f}_n(\theta) - f(\theta)| \right] \leq C w_n \sqrt{n},$$

where  $C$  is the same as in Lemma B.6, and  $\sup_{\theta \in \Theta} |\bar{f}_n(\theta) - f(\theta)|$  converges almost surely to zero.

*Proof.* We prove the two parts of the lemma separately. As in Lemma B.6, we assume all the quantities of interest to be measurable.

**First part of the lemma:**

Let us fix  $n > 0$ . It is easy to show that  $\bar{f}_n$  can be written as  $\bar{f}_n = \sum_{i=1}^n w_n^i f_i$  for some non-negative weights  $w_n^i$  with  $w_n^n = w_n$ . Let us also define the empirical cost  $F_i \triangleq \frac{1}{n-i+1} \sum_{j=i}^n f_j$ . According to (9), we have  $\mathbb{E} [\sup_{\theta \in \Theta} |F_i(\theta) - f(\theta)|] \leq \frac{C}{\sqrt{n-i+1}}$ . We now remark that

$$\bar{f}_n - f = \sum_{i=1}^n (w_n^i - w_n^{i-1})(n-i+1)(F_i - f),$$

where we have defined  $w_n^0 \triangleq 0$ . This relation can be proved by simple calculation. We obtain the first part by using the triangle inequality, and the fact that  $w_n^i \geq w_n^{i-1}$  for all  $i$ :

$$\begin{aligned}
\mathbb{E} \left[ \sup_{\theta \in \Theta} |\bar{f}_n(\theta) - f(\theta)| \right] &\leq \mathbb{E} \left[ \sum_{i=1}^n (w_n^i - w_n^{i-1})(n-i+1) \sup_{\theta \in \Theta} |F_i(\theta) - f(\theta)| \right] \\
&= \sum_{i=1}^n (w_n^i - w_n^{i-1})(n-i+1) \mathbb{E} \left[ \sup_{\theta \in \Theta} |F_i(\theta) - f(\theta)| \right] \\
&\leq \sum_{i=1}^n (w_n^i - w_n^{i-1}) C \sqrt{n-i+1} \\
&\leq \sqrt{n} C \sum_{i=1}^n (w_n^i - w_n^{i-1}) \\
&= C \sqrt{n} w_n.
\end{aligned}$$

This is unfortunately not sufficient to show that  $\mathbb{E} [\sup_{\theta \in \Theta} |\bar{f}_n(\theta) - f(\theta)|]$  converges to zero almost surely. We will show this fact by using Lemma A.6.

**Second part of the lemma:**

We call  $X_n = \sup_{\theta \in \Theta} |\bar{f}_n(\theta) - f(\theta)|$ . We have

$$\begin{aligned}
X_n - X_{n-1} &= \sup_{\theta \in \Theta} |(1-w_n)(\bar{f}_{n-1}(\theta) - f(\theta)) + w_n(f_n(\theta) - f(\theta))| - X_{n-1} \\
&\leq \sup_{\theta \in \Theta} w_n |f_n(\theta) - f(\theta)| - w_n X_{n-1} \leq 2M w_n
\end{aligned}$$

Let us denote by  $\theta_n^*$  a point in  $\Theta$  such that  $X_n = |\bar{f}_n(\theta_n^*) - f(\theta_n^*)|$ . We also have

$$\begin{aligned}
X_n - X_{n-1} &= \sup_{\theta \in \Theta} |(1-w_n)(\bar{f}_{n-1}(\theta) - f(\theta)) + w_n(f_n(\theta) - f(\theta))| - X_{n-1} \\
&\geq (1-w_n)X_{n-1} + w_n(f_n(\theta_{n-1}^*) - f(\theta_{n-1}^*)) - X_{n-1} \\
&\geq -w_n X_{n-1} + w_n(f_n(\theta_{n-1}^*) - f(\theta_{n-1}^*)) \\
&\geq -w_n 4M,
\end{aligned}$$

where we use again the fact that all functions  $f_n$ ,  $\bar{f}_n$  and  $f$  are bounded by  $M$ . Thus, we have shown that  $|X_n - X_{n-1}| \leq 4M w_n$ . Call  $a_n = w_n$  and  $b_n = w_n \sqrt{n}$ , then the conditions of Lemma A.6 are satisfied, and  $X_n$  converges almost surely to zero.  $\square$

Finally, the next lemma illustrates why the strong convexity of the surrogates is important.

**Lemma B.8 (Stability of the Estimates).**

When  $g_n$  is in  $\mathcal{S}_{L,\rho}(f, \theta_{n-1})$ ,

$$\|\theta_n - \theta_{n-1}\|_2 \leq \frac{2Rw_n}{\rho}.$$

*Proof.* Because the surrogates  $g_n$  are  $\rho$ -strongly convex, we have from Lemma A.4

$$\begin{aligned}
\frac{\rho}{2} \|\theta_n - \theta_{n-1}\|_2^2 &\leq \bar{g}_n(\theta_{n-1}) - \bar{g}_n(\theta_n) \\
&= w_n (g_n(\theta_{n-1}) - g_n(\theta_n)) + (1-w_n) (\bar{g}_{n-1}(\theta_{n-1}) - \bar{g}_{n-1}(\theta_n)) \\
&\leq w_n (g_n(\theta_{n-1}) - g_n(\theta_n)) \\
&\leq w_n (f_n(\theta_{n-1}) - f_n(\theta_n)) \\
&\leq R w_n \|\theta_n - \theta_{n-1}\|_2.
\end{aligned}$$

The second inequality comes from the fact that  $\theta_{n-1}$  is a minimizer of  $\bar{g}_{n-1}$ ; the third inequality is because  $g_n(\theta_{n-1}) = f_n(\theta_{n-1})$  and  $g_n \geq f_n$ . This is sufficient to conclude.  $\square$



## C Proofs of the Main Lemmas and Propositions

### C.1 Proof of Proposition 3.1

*Proof.* According to Lemma B.4, we have for all  $n \geq 1$ ,

$$w_n B_{n-1} \leq w_n f^* + L w_n A_{n-1} - L w_n \xi_{n-1} + w_n C_{n-1}.$$

By using the relations (7), this is equivalent to

$$B_{n-1} - B_n + w_n \mathbb{E}[f(\theta_{n-1})] \leq w_n f^* + L(A_{n-1} - A_n) + C_{n-1} - C_n + \frac{(R w_n)^2}{2L}.$$

By summing these inequalities between 1 and  $n$ , we obtain

$$B_0 - B_n + \sum_{k=1}^n w_k \mathbb{E}[f(\theta_{k-1})] \leq \left( \sum_{k=1}^n w_k \right) f^* + L A_0 - L A_n - C_n + \sum_{k=1}^n \frac{(R w_k)^2}{2L}.$$

Note that we also have

$$B_n \leq f^* + L A_n + C_n = L A_n + C_n + B_0 - L A_0 + L \xi_0.$$

Therefore, by combining the two previous inequalities,

$$\sum_{k=1}^n w_k \mathbb{E}[f(\theta_{k-1})] \leq \left( \sum_{k=1}^n w_k \right) f^* + L \xi_0 + \sum_{k=1}^n \frac{(R w_k)^2}{2L},$$

and by using Jensen's inequality,

$$\mathbb{E}[f(\bar{\theta}_{n-1}) - f^*] \leq \frac{L \xi_0 + \frac{R^2}{2L} \sum_{k=1}^n w_k^2}{\sum_{k=1}^n w_k}.$$

□

### C.2 Proof of Corollary 3.1

*Proof.*

We choose weights of the form  $w_n \triangleq \frac{\gamma}{\sqrt{n}}$ . Then, we have

$$\sum_{k=1}^n w_k^2 \leq \gamma^2 (1 + \log n),$$

by using the fact that  $\sum_{k=1}^n \frac{1}{k} \leq 1 + \log(n)$ . We also have for  $n \geq 2$ ,

$$\sum_{k=1}^n w_k \geq 2\gamma(\sqrt{n+1} - 1) \geq \gamma\sqrt{n},$$

where we use the fact that  $\sum_{k=1}^n \frac{1}{\sqrt{k}} \geq 2(\sqrt{n+1} - 1)$ , and the fact that  $2(\sqrt{n+1} - 1) \geq \sqrt{n}$  for all  $n \geq 2$ . Plugging this inequalities into (3) yields the desired result. □

### C.3 Proof of Proposition 3.2

*Proof.* We proceed in several steps, proving the convergence rates of several quantities of interest.

**Convergence rate of  $C_n$ :**

Let us show by induction that we have  $C_n \leq \frac{R^2}{\rho} w_n$  for all  $n \geq 1$ . This is obviously true for  $n = 1$

by definitions of  $w_1 = 1$  and  $C_1 = \frac{R^2}{2\rho}$ . Let us now assume that it is true for  $n - 1$ . We have

$$\begin{aligned}
C_n &= (1 - w_n)C_{n-1} + \frac{R^2}{2\rho}w_n^2 \\
&\leq \frac{R^2}{\rho}w_n \left( (1 - w_n)\frac{w_{n-1}}{w_n} + \frac{w_n}{2} \right) \\
&\leq \frac{R^2}{\rho}w_n \left( \frac{\beta(n-1)}{\beta n + 1} \frac{\beta n + 1}{\beta(n-1) + 1} + \frac{1}{\beta n + 1} \right) \\
&\leq \frac{R^2}{\rho}w_n \left( \frac{\beta(n-1)}{\beta(n-1) + 1} + \frac{1}{\beta(n-1) + 1} \right) \\
&= \frac{R^2}{\rho}w_n.
\end{aligned} \tag{10}$$

We conclude by induction that this is true for all  $n \geq 1$ .

**Convergence rate of  $A_n$ :**

From Lemma B.5 and B.4, we have for all  $n \geq 2$ ,

$$\mu A_{n-1} \leq L A_{n-1} - \rho \xi_{n-1} + C_{n-1}.$$

Multiplying this inequality by  $w_n$ ,

$$2\mu w_n A_{n-1} \leq \rho w_n (A_{n-1} - \xi_{n-1}) + w_n C_{n-1},$$

where the factor 2 comes from the fact that  $\rho = L + \mu$ . By using the definition of  $A_n$  in Eq. (7), we obtain the relation

$$A_n \leq \left(1 - \frac{2\mu w_n}{\rho}\right) A_{n-1} + \frac{w_n}{\rho} C_{n-1}.$$

Let us now show by induction that we have, for all  $n \geq 1$ , the convergence rate  $A_n \leq \delta w_n$ , where  $\delta \triangleq \max\left(\frac{R^2}{\rho\mu}, \xi_0\right)$ . For  $n = 1$ , we have that  $w_1 = 1$ , and thus  $A_1 = \xi_0 \leq \delta$ . Assume now that we have  $A_{n-1} \leq \delta w_{n-1}$  for some  $n \geq 1$ . Then, by using the convergence rate (10) and the induction hypothesis,

$$\begin{aligned}
A_n &\leq \delta w_n \left( \left(1 - \frac{2\mu w_n}{\rho}\right) \frac{w_{n-1}}{w_n} + \frac{R^2 w_{n-1}}{\rho^2 \delta} \right) \\
&\leq \delta w_n \left( \left(1 - \frac{2\mu w_n}{\rho}\right) \frac{w_{n-1}}{w_n} + \mu \frac{w_{n-1}}{\rho} \right) \\
&\leq \delta w_n \left( \frac{\beta n + 1 - \frac{2\mu(1+\beta)}{\rho}}{\beta n + 1} \frac{\beta n + 1}{\beta(n-1) + 1} + \frac{\frac{\mu(1+\beta)}{\rho}}{\beta(n-1) + 1} \right) \\
&= \delta w_n \left( \frac{\beta n + 1 - \frac{\mu(1+\beta)}{\rho}}{\beta(n-1) + 1} \right) \\
&\leq \delta w_n.
\end{aligned}$$

The last inequality uses the fact that  $\frac{\mu(1+\beta)}{\rho} \geq \beta$  because  $\beta \leq \frac{\mu}{L}$ . We conclude by induction that  $A_n \leq \delta w_n$  for all  $n \geq 1$ .

**Convergence rate of  $\mathbb{E}[f(\hat{\theta}_n) - f^*] + \rho \xi_n$ :**

We use again Lemma B.4:

$$B_n - f^* + \rho \xi_n \leq L A_n + C_n,$$

and we consider two possible cases

- If  $\frac{R^2}{\rho\mu} \geq \xi_0$ , then

$$\begin{aligned} B_n - f^* + \rho\xi_n &\leq \frac{R^2}{\rho} \left(1 + \frac{L}{\mu}\right) w_n \\ &= \frac{R^2}{\mu} w_n \\ &\leq \frac{2R^2}{\mu(\beta n + 1)}, \end{aligned}$$

where we simply use the convergence rates of  $A_n$  and  $C_n$  computed before.

- If instead  $\frac{R^2}{\rho\mu} < \xi_0$ , then

$$\begin{aligned} B_n - f^* + \rho\xi_n &\leq \left(\frac{R^2}{\rho} + L\xi_0\right) w_n \\ &\leq \rho\xi_0 w_n \\ &\leq \frac{2\rho\xi_0}{\beta n + 1}. \end{aligned}$$

It is then easy to prove that  $\mathbb{E}[f(\hat{\theta}_n) - f^*] \leq B_n$  by using Jensen's inequality, which allows us to conclude. □

#### C.4 Proof of Proposition 3.3

*Proof.* We generalize the proof of convergence for online matrix factorization of [19]. The proof exploits Theorem A.1 about the convergence of quasi-martingales [33], similarly as [3] for proving the convergence of the stochastic gradient descent algorithm for non-convex functions.

**Almost sure convergence of  $(\bar{g}_n(\theta_n))_{n \geq 1}$ :**

The first step consists of applying a convergence theorem for the sequence  $(\bar{g}_n(\theta_n))_{n \geq 1}$  by bounding its positive expected variations. Define  $Y_n \triangleq \bar{g}_n(\theta_n)$ . For  $n \geq 2$ , we have

$$\begin{aligned} Y_n - Y_{n-1} &= \bar{g}_n(\theta_n) - \bar{g}_n(\theta_{n-1}) + \bar{g}_n(\theta_{n-1}) - \bar{g}_{n-1}(\theta_{n-1}) \\ &= (\bar{g}_n(\theta_n) - \bar{g}_n(\theta_{n-1})) + w_n(g_n(\theta_{n-1}) - \bar{g}_{n-1}(\theta_{n-1})) \\ &= (\bar{g}_n(\theta_n) - \bar{g}_n(\theta_{n-1})) + w_n(\bar{f}_{n-1}(\theta_{n-1}) - \bar{g}_{n-1}(\theta_{n-1})) + w_n(g_n(\theta_{n-1}) - \bar{f}_{n-1}(\theta_{n-1})) \\ &= (\bar{g}_n(\theta_n) - \bar{g}_n(\theta_{n-1})) + w_n(\bar{f}_{n-1}(\theta_{n-1}) - \bar{g}_{n-1}(\theta_{n-1})) + w_n(f_n(\theta_{n-1}) - \bar{f}_{n-1}(\theta_{n-1})) \\ &\leq w_n(f_n(\theta_{n-1}) - \bar{f}_{n-1}(\theta_{n-1})). \end{aligned} \tag{11}$$

The final inequality comes from the inequality  $\bar{g}_n \geq \bar{f}_n$ , which is easy to show by induction starting from  $n = 1$  since  $w_1 = 1$ . It follows,

$$\begin{aligned} \mathbb{E}[\bar{g}_n(\theta_n) - \bar{g}_{n-1}(\theta_{n-1}) | \mathcal{F}_{n-1}] &\leq w_n \mathbb{E}[f_n(\theta_{n-1}) - \bar{f}_{n-1}(\theta_{n-1}) | \mathcal{F}_{n-1}] \\ &= w_n(f(\theta_{n-1}) - \bar{f}_{n-1}(\theta_{n-1})) \\ &\leq w_n \sup_{\theta \in \Theta} |f(\theta) - \bar{f}_{n-1}(\theta)|, \end{aligned}$$

where  $\mathcal{F}_{n-1}$  is the filtration representing the past information before time  $n$ . Call now

$$\delta_n \triangleq \begin{cases} 1 & \text{if } \mathbb{E}[\bar{g}_n(\theta_n) - \bar{g}_{n-1}(\theta_{n-1}) | \mathcal{F}_{n-1}] > 0 \\ 0 & \text{otherwise.} \end{cases}$$

Then, the series below with non-negative summands converges:

$$\begin{aligned} \sum_{n=1}^{\infty} \mathbb{E}[\delta_n(\bar{g}_n(\theta_n) - \bar{g}_{n-1}(\theta_{n-1}))] &= \sum_{n=1}^{\infty} \mathbb{E}[\delta_n \mathbb{E}[(\bar{g}_n(\theta_n) - \bar{g}_{n-1}(\theta_{n-1})) | \mathcal{F}_{n-1}]] \\ &\leq \sum_{n=1}^{\infty} \mathbb{E} \left[ w_n \sup_{\theta \in \Theta} |f(\theta) - \bar{f}_{n-1}(\theta)| \right] \\ &\leq \sum_{n=1}^{\infty} C w_n^2 \sqrt{n} < +\infty, \end{aligned}$$

The second inequality comes from Lemma B.7. Since in addition  $\bar{g}_n$  is bounded below by some constant independent of  $n$ , we can apply Theorem A.1. This theorem tells us that  $(\bar{g}_n(\theta_n))_{n \geq 1}$  converges almost surely to an integrable random variable  $g^*$  and that  $\sum_{n=1}^{\infty} \mathbb{E}[|\mathbb{E}[\bar{g}_n(\theta_n) - \bar{g}_{n-1}(\theta_{n-1}) | \mathcal{F}_{n-1}]|]$  converges almost surely.

**Almost sure convergence of  $(\bar{f}_n(\theta_n))_{n \geq 1}$ :**

We will show by using Lemma A.5 that the non-positive term  $\bar{f}_n(\theta_n) - \bar{g}_n(\theta_n)$  almost surely converges to zero, and thus  $(\bar{f}_n(\theta_n))_{n \geq 1}$  is also converging almost surely to  $g^*$ .

We observe that

$$\sum_{n=1}^{\infty} \mathbb{E}[|\mathbb{E}[\bar{g}_n(\theta_n) - \bar{g}_{n-1}(\theta_{n-1}) | \mathcal{F}_{n-1}]|] = \mathbb{E} \left[ \sum_{n=1}^{\infty} |\mathbb{E}[\bar{g}_n(\theta_n) - \bar{g}_{n-1}(\theta_{n-1}) | \mathcal{F}_{n-1}]| \right] < +\infty.$$

Thus, the series  $\sum_{n=1}^{\infty} |\mathbb{E}[\bar{g}_n(\theta_n) - \bar{g}_{n-1}(\theta_{n-1}) | \mathcal{F}_{n-1}]|$  is absolutely convergent with probability one, and the series  $\sum_{n=1}^{\infty} \mathbb{E}[\bar{g}_n(\theta_n) - \bar{g}_{n-1}(\theta_{n-1}) | \mathcal{F}_{n-1}]$  is also almost surely convergent.

We also remark that, using Lemma B.7,

$$\mathbb{E} \left[ \sum_{n=1}^{+\infty} w_n |f(\theta_{n-1}) - \bar{f}_{n-1}(\theta_{n-1})| \right] \leq C \sum_{n=1}^{+\infty} w_n^2 \sqrt{n} < +\infty,$$

and thus  $w_n(f(\theta_{n-1}) - \bar{f}_{n-1}(\theta_{n-1}))$  is the summand of an absolutely convergent series with probability one.

Taking the expectation of Eq. (11) conditioned on  $\mathcal{F}_{n-1}$ , it remains that the non-positive term  $w_n(\bar{f}_{n-1}(\theta_{n-1}) - \bar{g}_{n-1}(\theta_{n-1}))$  is also necessarily the summand of an almost surely convergent series, since all other terms in the equation are summands of almost surely converging sums. This is not sufficient to immediately conclude that  $\bar{f}_n(\theta_n) - \bar{g}_n(\theta_n)$  converges to zero almost surely, and thus we will use Lemma A.5. We have that  $\sum_{n=1}^{+\infty} w_n$  diverges, that  $\sum_{n=1}^{+\infty} w_n(\bar{g}_{n-1}(\theta_{n-1}) - \bar{f}_{n-1}(\theta_{n-1}))$  converges almost surely. Define  $X_n \triangleq (\bar{g}_{n-1}(\theta_{n-1}) - \bar{f}_{n-1}(\theta_{n-1}))$ . By definition of the surrogate functions, the differences  $h_n \triangleq g_n - f_n$  are differentiable and their gradients are  $L$ -Lipschitz continuous. Since in addition  $\Theta$  is compact and  $\nabla h_n(\theta_{n-1}) = 0$ ,  $\nabla h_n$  is bounded by some constant  $R'$  independent of  $n$ , and the function  $h_n$  is  $R'$ -Lipschitz. This is therefore also the case for  $\bar{h}_n = \bar{g}_n - \bar{f}_n$ .

$$\begin{aligned} |X_{n+1} - X_n| &= |\bar{h}_n(\theta_n) - \bar{h}_{n-1}(\theta_{n-1})| \\ &\leq |\bar{h}_n(\theta_n) - \bar{h}_n(\theta_{n-1})| + |\bar{h}_n(\theta_{n-1}) - \bar{h}_{n-1}(\theta_{n-1})| \\ &\leq R' \|\theta_n - \theta_{n-1}\|_2 + |\bar{h}_n(\theta_{n-1}) - \bar{h}_{n-1}(\theta_{n-1})| \\ &\leq \frac{2RR'}{\rho} w_n + w_n |h_n(\theta_{n-1}) - \bar{h}_{n-1}(\theta_{n-1})| \\ &= \frac{2RR'}{\rho} w_n + w_n |\bar{h}_{n-1}(\theta_{n-1})| \\ &\leq O(w_n). \end{aligned}$$

The second inequality uses the fact that  $\bar{h}_n$  is  $R'$ -Lipschitz; The second inequality uses Lemma B.8; the last equality uses the fact that the functions  $h_n$  are also bounded by some constant independent of  $n$  (using the fact that  $\nabla h_n$  is uniformly bounded). We can now apply Lemma A.5, and  $X_n$  converges to zero with probability one. Thus,  $(\bar{f}_n(\theta_n))_{n \geq 1}$  converges almost surely to  $g^*$ .

**Almost sure convergence of  $(f(\theta_n))_{n \geq 1}$ :**

Since  $(\bar{f}_n(\theta_n))_{n \geq 1}$  converges almost surely, we simply use Lemma A.6, which tells us that  $\bar{f}_n$  converges uniformly to  $f$ . Then,  $(f(\theta_n))_{n \geq 1}$  converges almost surely to  $g^*$ .

**Asymptotic Stationary Point Condition:**

Let us call  $\bar{h}_n \triangleq \bar{g}_n - \bar{f}_n$ , which can be shown to be differentiable with a  $L$ -Lipschitz gradient by definition of the surrogate  $g_n$ . For all  $\theta$  in  $\Theta$ ,

$$\nabla \bar{f}_n(\theta_n, \theta - \theta_n) = \nabla \bar{g}_n(\theta_n, \theta - \theta_n) - \nabla \bar{h}_n(\theta_n)^\top (\theta - \theta_n).$$

Since  $\theta_n$  is the minimizer of  $\bar{g}_n$ , we have  $\nabla \bar{g}_n(\theta_n, \theta - \theta_n) \geq 0$ .

Since  $\bar{h}_n$  is differentiable and its gradient is  $L$ -Lipschitz continuous, we can apply Lemma A.1 to  $\theta = \theta_n$  and  $\theta' = \theta_n - \frac{1}{L} \nabla \bar{h}_n(\theta_n)$ , which gives  $\bar{h}_n(\theta') \leq \bar{h}_n(\theta_n) - \frac{1}{2L} \|\nabla \bar{h}_n(\theta_n)\|_2^2$ . Since we have shown that  $\bar{h}_n(\theta_n) = \bar{g}_n(\theta_n) - f(\theta_n)$  converges to zero and  $\bar{h}_n(\theta') \geq 0$ , we have that  $\|\nabla \bar{h}_n(\theta_n)\|_2$  converges to zero. Thus,

$$\inf_{\theta \in \Theta} \frac{\nabla \bar{f}_n(\theta_n, \theta - \theta_n)}{\|\theta - \theta_n\|_2} \geq -\|\nabla \bar{h}_n(\theta_n)\|_2 \xrightarrow{n \rightarrow +\infty} 0 \text{ a.s.}$$

□

**C.5 Proof of Proposition 3.4**

*Proof.* Since  $\Theta$  is compact according to assumption (C), the sequence  $(\theta_n)_{n \geq 1}$  admits limit points. Let us consider a converging subsequence  $(n_k)_{k \geq 1}$  to a limit point  $\theta_\infty$  in  $\Theta$ . In this converging subsequence, we can also find a subsequence  $(n_{k'})_{k' \geq 1}$  such that  $\kappa_{n_{k'}}$  converges to a point  $\kappa_\infty$  in  $\mathcal{K}$  (which is compact). For the sake of simplicity, and without loss of generality, we remove the indices  $k$  and  $k'$  from the notation and assume that  $\theta_n$  converges to  $\theta_\infty$ , while  $\kappa_n$  converges to  $\kappa_\infty$ . It is then easy to see that the functions  $\bar{g}_n$  converge uniformly to  $\bar{g}_\infty \triangleq g_{\kappa_\infty}$ , given the assumptions made in the proposition.

Defining  $\bar{h}_\infty \triangleq \bar{g}_\infty - f$ , we have for all  $\theta$  in  $\Theta$ :

$$\nabla f(\theta_\infty, \theta - \theta_\infty) = \nabla \bar{g}_\infty(\theta_\infty, \theta - \theta_\infty) - \nabla \bar{h}_\infty(\theta_\infty, \theta - \theta_\infty).$$

To prove the proposition, we will first show that  $\nabla \bar{g}_\infty(\theta_\infty, \theta - \theta_\infty) \geq 0$  and then that  $\nabla \bar{h}_\infty(\theta_\infty, \theta - \theta_\infty) = 0$ .

**Proof of  $\nabla \bar{g}_\infty(\theta_\infty, \theta - \theta_\infty) \geq 0$ :**

It is sufficient to show that  $\theta_\infty$  is a minimizer of  $\bar{g}_\infty$ . This is straightforward, by taking the limit when  $n$  goes to infinity of

$$\bar{g}_n(\theta) \geq \bar{g}_n(\theta_n),$$

where we use the uniform convergence of  $\bar{g}_n$ .

**Proof of  $\nabla \bar{h}_\infty(\theta_\infty, \theta - \theta_\infty) = 0$ :**

Since both  $\bar{f}_n$  and  $\bar{g}_n$  converges uniformly (according to Lemma B.7 for  $\bar{f}_n$ ), we have that  $\bar{h}_n$  converges uniformly to  $\bar{h}_\infty$ . Since  $\bar{h}_n$  is differentiable with a  $L$ -Lipschitz gradient, we have for all vector  $\mathbf{z}$  in  $\mathbb{R}^p$ ,

$$\bar{h}_n(\theta_n + \mathbf{z}) = \bar{h}_n(\theta_n) + \nabla \bar{h}_n(\theta_n)^\top \mathbf{z} + O(\|\mathbf{z}\|_2^2),$$

where the constant in  $O$  is independent of  $n$ . By taking the limit when  $n$  goes to infinity, it remains

$$\bar{h}_\infty(\theta_\infty + \mathbf{z}) = \bar{h}_\infty(\theta_\infty) + O(\|\mathbf{z}\|_2^2),$$

since we have shown in the proof of Proposition 3.3 that  $\|\nabla \bar{h}_n(\theta_n)\|_2$  converges to zero. Since  $\bar{h}_\infty$  admits a first order extension around  $\theta_\infty$  it is differentiable at this point and furthermore,  $\nabla \bar{h}_\infty(\theta_\infty) = 0$ . This is sufficient to conclude. □

**C.6 Proof of Proposition 3.5**

*Proof.* First we notice that

- $g_n \geq f_n$ ;

- $g_n(\theta_{n-1}) = f_n(\theta_{n-1})$ ;
- $g_n$  is  $\rho_1$ -strongly convex since  $\theta \mapsto g_{k,n}(\gamma_k(\theta))$  can be shown to be convex, following elementary composition rules for convex functions (see [32], Section 3.2.4).

Thus, the only property missing is the smoothness of the approximation error  $h_n \triangleq g_n - f_n$ . Rather than writing again a full proof, we now simply review the different places where this property is used, and which modifications should be made to the proofs of Propositions 3.3 and 3.4.

In the second step of this proof, we require the functions  $\bar{h}_n$  to be uniformly Lipschitz and uniformly bounded. It is easy to check that it is still the case with the assumptions we made in Proposition 3.5.

The last step about the asymptotic point condition is however more problematic, where we cannot show anymore that the quantity  $\nabla \bar{h}_n(\theta_n)$  converges to zero (since  $\bar{h}_n$  is not differentiable anymore). Instead, we need to show that the directional derivative  $\frac{\nabla \bar{h}_n(\theta_n, \theta - \theta_n)}{\|\theta - \theta_n\|}$  uniformly converges to zero on  $\Theta$ .

We will show the result for  $K = 1$ ; it will be easy to extend it to any arbitrary  $K > 2$ . We remark that

$$\nabla \bar{h}_n(\theta_n, \theta - \theta_n) = \nabla \bar{h}_{0,n}(\theta_n)^\top (\theta - \theta_n) + \lim_{t \rightarrow 0^+} \frac{\bar{h}_{1,n}(\gamma_1(\theta_n + t(\theta - \theta_n))) - \bar{h}_{1,n}(\gamma_1(\theta_n))}{t},$$

where  $\bar{h}_{0,n}$  and  $\bar{h}_{1,n}$  are defined similarly as  $\bar{h}_n$  for the functions  $h_{0,n} \triangleq g_{0,n} - f_{0,n}$  and  $h_{1,n} \triangleq g_{1,n} - f_{1,n}$  respectively. Since  $\bar{h}_n(\theta_n)$  is shown to converge to zero, we have that the non-negative quantities  $\bar{h}_{0,n}(\theta_n)$  and  $\bar{h}_{1,n}(\gamma_1(\theta_n))$  converge to zero as well. Since  $\bar{h}_{0,n}$  and  $\bar{h}_{1,n}$  are differentiable and their gradients are Lipschitz, we use similar arguments as in the proof of Proposition 3.3, and we have that  $\nabla \bar{h}_{0,n}(\theta_n)$  and  $\bar{h}'_{1,n}(\gamma_1(\theta_n))$  converge to zero (where  $\bar{h}'_{1,n}$  is the derivative of  $\bar{h}_{1,n}$ ). Concerning the second term, we can make the following Taylor expansion for  $\bar{h}_{1,n}$ :

$$\bar{h}_{1,n}(\gamma_1(\theta_n + \mathbf{z})) = \bar{h}_{1,n}(\gamma_1(\theta_n)) + \bar{h}'_{1,n}(\gamma_1(\theta_n))(\gamma_1(\theta_n + \mathbf{z}) - \gamma_1(\theta_n)) + O((\gamma_1(\theta_n + \mathbf{z}) - \gamma_1(\theta_n))^2),$$

where the constant in the  $O$  notation is independent of  $\theta_n$  and  $\mathbf{z}$  (since the derivative is  $L_1$ -Lipschitz). Plugging  $\mathbf{z} \triangleq t(\theta - \theta_n)$  in this last equation, and using the Lipschitz property of  $\gamma_1$ , we have

$$\lim_{t \rightarrow 0^+} \left| \frac{\bar{h}_{1,n}(\gamma_1(\theta_n + t(\theta - \theta_n))) - \bar{h}_{1,n}(\gamma_1(\theta_n))}{t} \right| \leq |\bar{h}'_{1,n}(\gamma_1(\theta_n))| \|\theta - \theta_n\|.$$

Since  $\bar{h}'_{1,n}(\gamma_1(\theta_n))$  converges to zero, we can conclude the proof of the modified Proposition 3.3.

The proof of Proposition 3.4 can be modified with similar arguments.  $\square$

## D Additional Experimental Results

We present in Figures 4 and 5 some additional experimental comparisons, which complement the ones of Section 4.1. Figures 6 and 7 present additional plots from the experiment of Section 4.2. Finally, we present three dictionaries corresponding to the experiment of Section 4.3 in Figures 8, 9 and 10.

## Supplementary References

- [31] D.P. Bertsekas. *Nonlinear programming*. Athena Scientific Belmont, 1999. 2nd edition.
- [32] S.P. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [33] D. L. Fisk. Quasi-martingales. *T. Am. Math. Soc.*, 120(3):359–388, 1965.
- [34] M. Métivier. *Semi-martingales*. Walter de Gruyter, 1983.
- [35] Y. Nesterov. *Introductory lectures on convex optimization*. Kluwer Academic Publishers, 2004.
- [36] Y. Nesterov and J.-P. Vial. Confidence level solutions for stochastic programming. *Automatica*, 44(6):1559–1568, 2008.
- [37] J. Nocedal and S.J. Wright. *Numerical optimization*. Springer Verlag, 2006. 2nd edition.

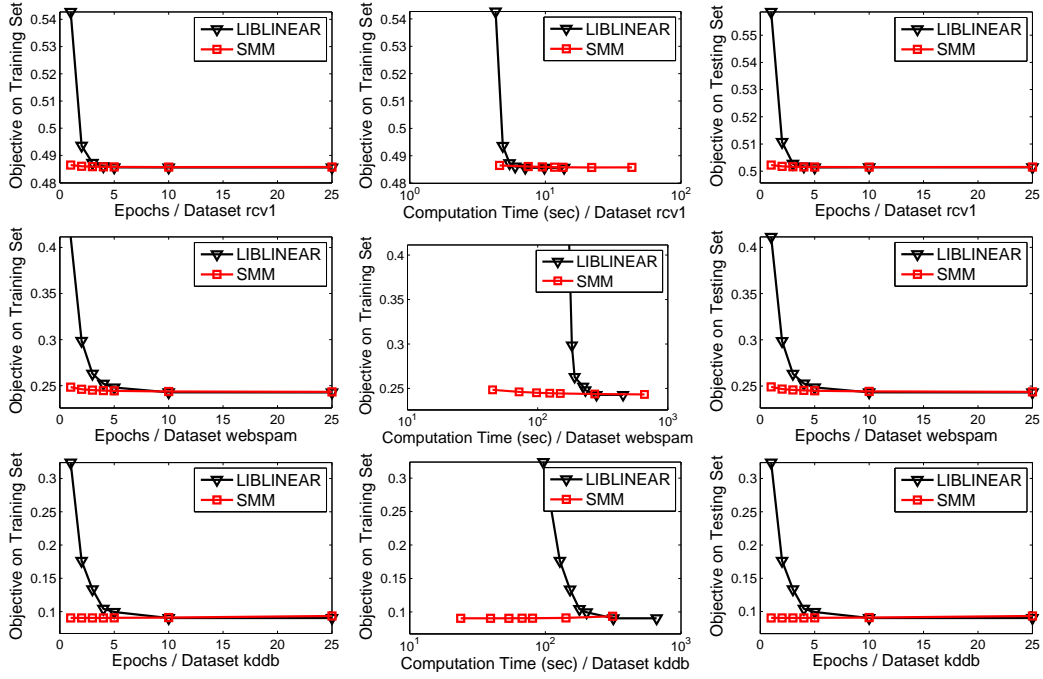


Figure 4: Comparison between LIBLINEAR and SMM in the high regularization regime for  $\ell_1$ -logistic regression.

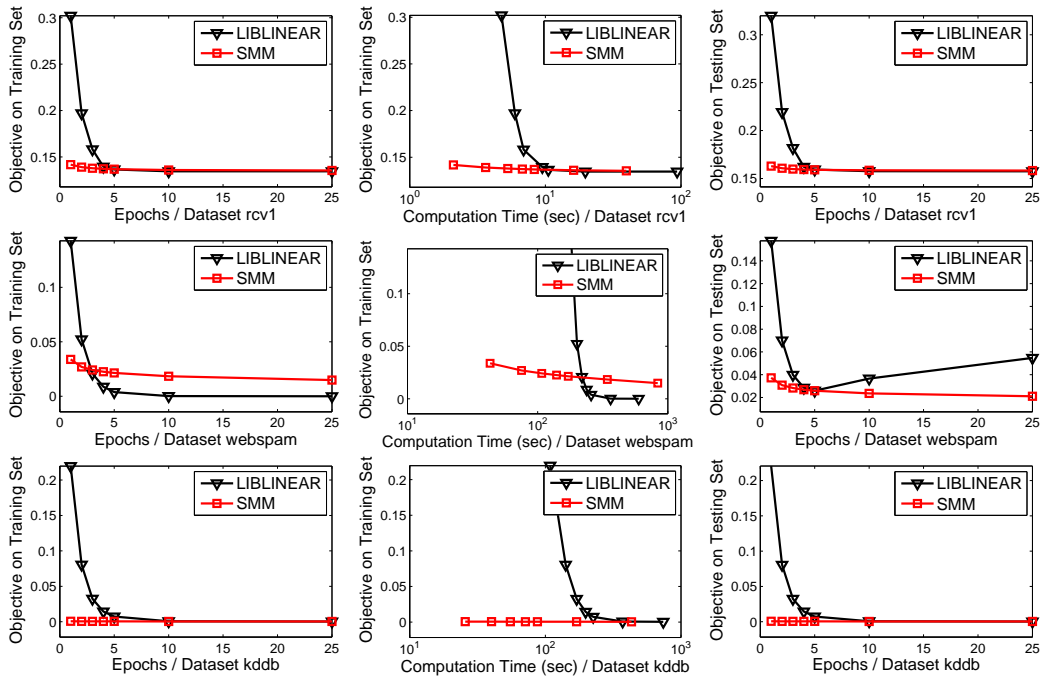


Figure 5: Comparison between LIBLINEAR and SMM in the low regularization regime for  $\ell_1$ -logistic regression.

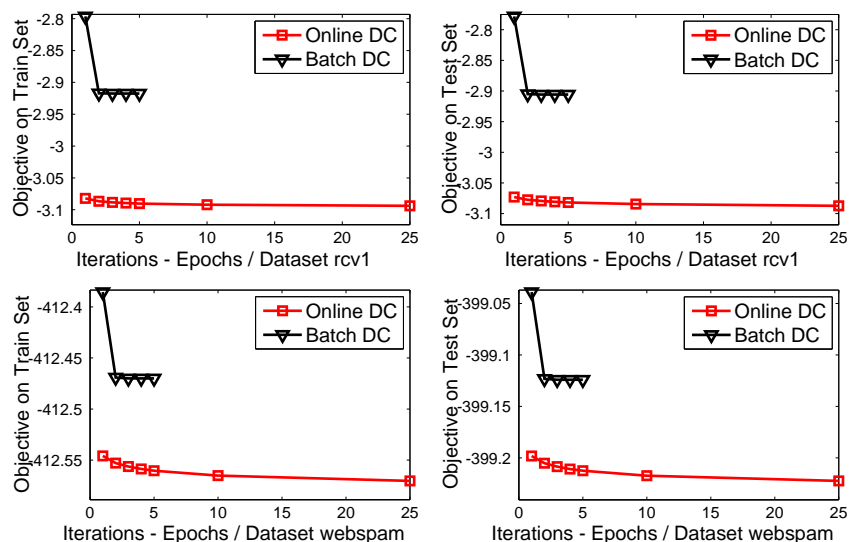


Figure 6: Comparison between batch and online DC programming, with high regularization for the datasets rcv1 and webspam. Note that each iteration in the batch setting can perform several epochs.

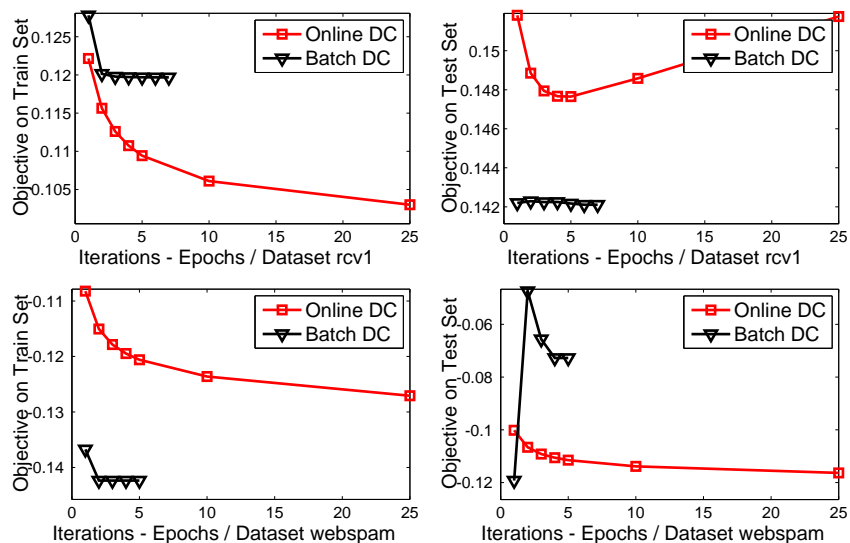


Figure 7: Comparison between batch and online DC programming, with low regularization for the datasets rcv1 and webspam. Note that each iteration in the batch setting can perform several epochs.



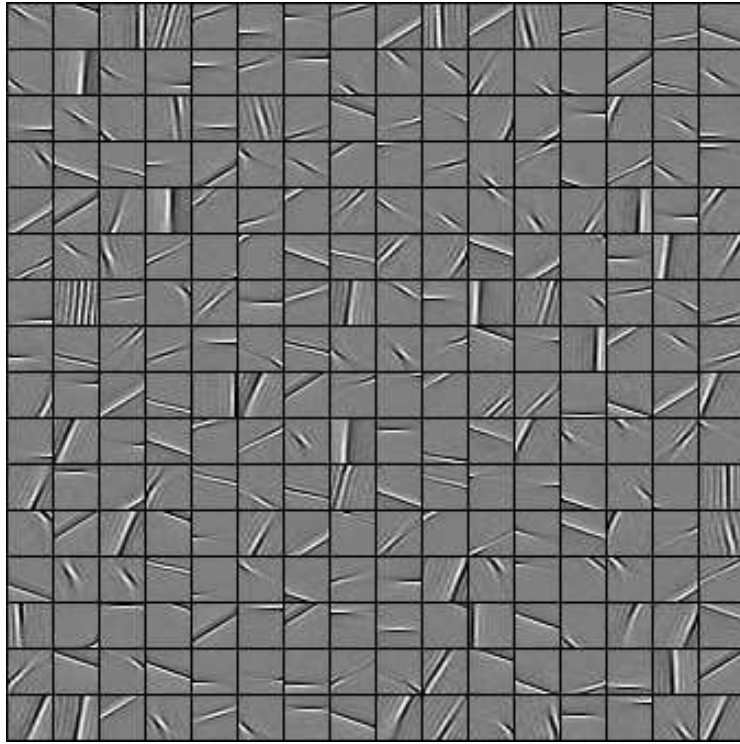


Figure 8: Dictionary obtained using the toolbox SPAMS [19].

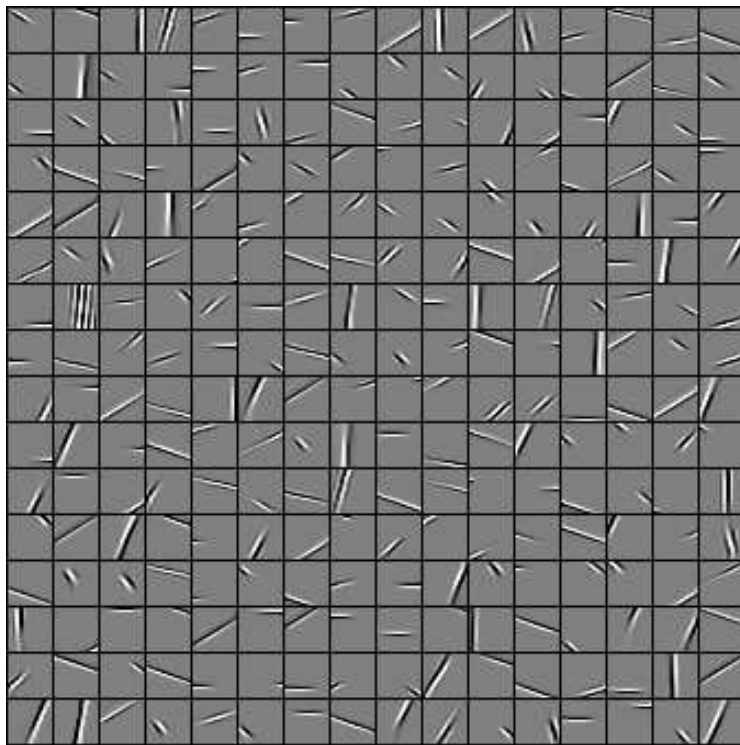


Figure 9: Sparse dictionary obtained by our approach, using the dictionary of Figure 8 as an initialization.

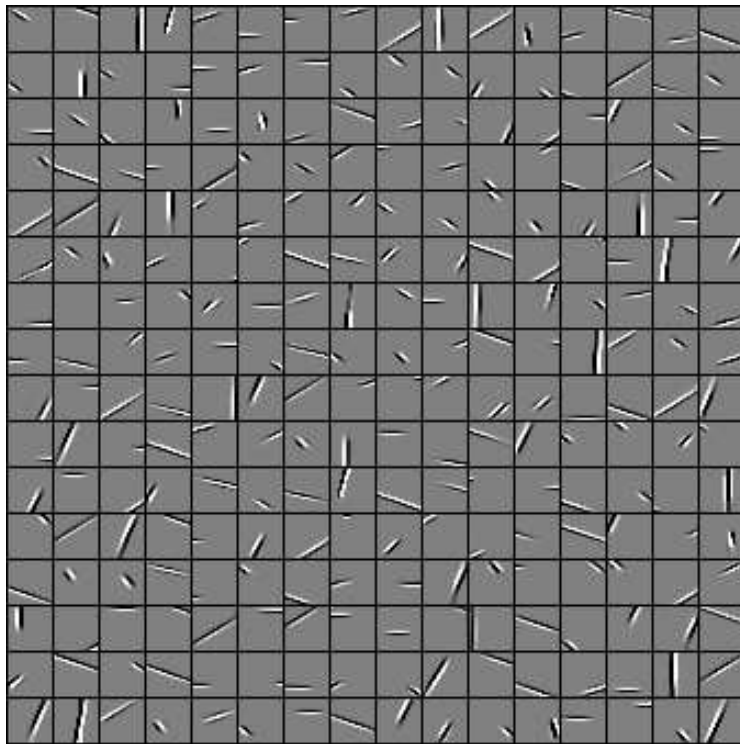


Figure 10: Sparse dictionary obtained by our approach, using the dictionary of Figure 8 as an initialization, and with a higher regularization parameter than in Figure 9.