



HAL
open science

Incremental Without Replacement Sampling in Nonconvex Optimization

Edouard Pauwels

► **To cite this version:**

Edouard Pauwels. Incremental Without Replacement Sampling in Nonconvex Optimization. Journal of Optimization Theory and Applications, 2021, 10.1007/s10957-021-01883-2 . hal-02896102v4

HAL Id: hal-02896102

<https://hal.science/hal-02896102v4>

Submitted on 26 Dec 2022

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.

Incremental Without Replacement Sampling in Nonconvex Optimization

Edouard Pauwels*

December 26, 2022

Abstract

Minibatch decomposition methods for empirical risk minimization are commonly analysed in a stochastic approximation setting, also known as sampling with replacement. On the other hands modern implementations of such techniques are incremental: they rely on sampling without replacement, for which available analysis are much scarcer. We provide convergence guaranties for the latter variant by analysing a versatile incremental gradient scheme. For this scheme, we consider constant, decreasing or adaptive step sizes. In the smooth setting we obtain explicit complexity estimates in terms of epoch counter. In the nonsmooth setting we prove that the sequence is attracted by solutions of optimality conditions of the problem.

Communicated by Gabriel Peyré

Keywords. Without Replacement Sampling, Incremental Methods, Nonconvex Optimization, First order Methods, Stochastic Gradient, Adaptive Methods, Backpropagation, Deep Learning

1 Introduction

1.1 Context and motivation

Training of modern learning architectures is mostly achieved by empirical risk minimization, relying on minibatch decomposition first order methods [20, 35]. The goal is to solve optimization problems of the form

$$F^* = \inf_{x \in \mathbb{R}^p} F(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \quad (1)$$

where $f_i: \mathbb{R}^p \rightarrow \mathbb{R}$ are Lipschitz functions and the infimum is finite. In this context minibatching takes advantage of redundancy in large sums and perform steps which only

*IRIT, Université de Toulouse, CNRS. Toulouse, France.

rely on partial sums [18]. The most widely studied variant is the Stochastic Gradient algorithm (also known as SGD), each step consists in sampling with replacement in $\{1, \dots, n\}$, and moving in the direction of the gradient of F corrupted by centered noise inherent to subsampling. This allows to study such algorithms in the broader context of stochastic approximation, initiated by Robins and Monro [51] with many subsequent works [37, 6, 33, 17, 42, 23, 14].

On the other hand most widely used implementations of such learning strategies for deep network [1, 46] rely on sampling without replacement, an epoch being the result of a single path during which first order information for each f_i is computed exactly once. Although very close to stochastic approximation, this strategy does not satisfy the “gradient plus centered noise” hypothesis. Therefore all existing theoretical guaranties relying on stochastic approximation arguments do not hold true for many practical implementation of learning algorithm. The purpose of this work is to provide convergence guaranties for such “without replacement minibatch strategies” for problem (1), also known as incremental methods [10]. The main strategy is to view the algorithmic steps through the lenses of perturbed gradient iterations, see for example [9].

1.2 Problem setting

We consider problem (1) and assume that for each f_i is Lipschitz and that we have access to an oracle which provides a search direction $d_i: \mathbb{R}^p \mapsto \mathbb{R}^p$, $i = 1, \dots, n$. We consider two settings

Smooth setting: f_i are C^1 with Lipschitz gradient, in which case we set $d_i = \nabla f_i$, $i = 1, \dots, n$.

Nonsmooth setting: f_i are path differentiable. Such functions constitute a subclass of Lipschitz functions which enjoy some of the nice properties of nonsmooth convex functions [14] in particular, an operational chain rule [23]. In this case, d_i is a selection in a conservative field for f_i . Examples of such objects include convex subgradients if f_i is convex, the Clarke subgradient, which is an extension to the nonconvex setting, as well as the output of automatic differentiation applied to a nonsmooth program, see [14] for more details.

We consider a class of descent methods described by Algorithm 1. Notice that there is no randomness specified in the algorithm, all our results are worst case and hold deterministically. Algorithm 1 allows to model:

Gradient descent: Set $\hat{z}_{K,i-1} = z_{K,0} = x_K$, for all $K \in \mathbb{N}$ and $i = 1 \dots, n$.

Incremental algorithms: Set $\hat{z}_{K,i-1} = z_{K,i-1}$, for all $K \in \mathbb{N}$ and $i = 1 \dots, n$.

Random permutations: Although not explicitly stated in the algorithm, all our proof argument hold independantly of the order of query of the indices $i = 1, \dots, n$ for each epoch. Hence our results actually hold deterministically for “random shuffling” or, “without replacement sampling” strategies.

Mini-batching: Set $\hat{z}_{K,i-1} = \hat{z}_{K,i-2} = z_{K,i-2}$, which results in computation of gradients of f_i and f_{i-1} at the same point $z_{K,i-2}$.

Algorithm 1: Without replacement descent algorithm

Program data: For $i = 1 \dots, n$, f_i , corresponding search direction oracle d_i .

Input: $x_0 \in \mathbb{R}^p$.

1: Decreasing steps:	1: Adaptive steps:
2: $(\alpha_{K,i})_{K \in \mathbb{N}, i \in \{1, \dots, n\}}$.	2: $v_0 = \delta > 0, \beta > 0$.
3: for $K \in \mathbb{N}$ do	3: for $K \in \mathbb{N}$ do
4: Set: $z_{K,0} = x_K$	4: Set: $z_{K,0} = x_K, v_{K,0} = v_K$
5: for $i = 1, \dots, n$ do	5: for $i = 1, \dots, n$ do
6: $\hat{z}_{K,i-1} \in \text{conv} \left((z_{K,j})_{j=0}^{i-1} \right)$.	6: $\hat{z}_{K,i-1} \in \text{conv} \left((z_{K,j})_{j=0}^{i-1} \right)$.
7: $z_{K,i} = z_{K,i-1} - \alpha_{K,i} d_i(\hat{z}_{K,i-1})$	7: $v_{K,i} = v_{K,i-1} + \beta \ d_i(\hat{z}_{K,i-1})\ ^2$
8: end for	8: $\alpha_{K,i} = v_{K,i}^{-1/3}$
9: Set: $x_{K+1} = z_{K,n}$	9: $z_{K,i} = z_{K,i-1} - \alpha_{K,i} d_i(\hat{z}_{K,i-1})$
10: end for	10: end for
	11: Set: $x_{K+1} = z_{K,n}, v_{K+1} = v_{K,n}$.
	12: end for

Asynchronous computation in a parameter server setting: Consider that $(z_{K,i})_{K \in \mathbb{N}, i=1 \dots n}$ is stored on a server, accessed by workers which compute $d_i(z_{K,i-1})$. Due to communication and computation delays, d_i may be evaluated using an outdated estimate of z , called \hat{z} . In Algorithm 1, asynchronicity and delays between workers may be arbitrary within each epoch. However, we enforce that the whole system waits for all workers to communicate results before starting a new epoch, a form of partial synchronization.

1.3 Contributions

We propose a detailed convergence analysis of Algorithm 1 in a nonconvex setting. Our analysis is worst case and our results hold deterministically. When each f_i is smooth with L_i -Lipschitz gradient, setting $L = \frac{1}{n} \sum_{i=1}^n L_i$, we obtain the following estimates on the squared norm of the gradient of F in terms of number of epochs K (all constants are explicit).

- Decreasing step size without knowledge of L : $O\left(\frac{1}{\sqrt{K}}\right)$.
- Decreasing step size with knowledge of L : $O\left(\frac{1}{K^{2/3}}\right)$.
- Adaptive step size without knowledge of L : $O\left(\frac{1}{K^{2/3}}\right)$.

For general nonsmooth objectives, convergence rate do not exist for the simplest subgradient oracle, see for example [57] with an attempt for more complex oracles. We prove that the sequence $(x_K)_{K \in \mathbb{N}}$ is attracted by subsets of \mathbb{R}^p which are solutions to optimality condition related to problem (1), for both step size strategies.

1.4 Relation to existing literature

Incremental gradient was introduced by Bertsekas in the late 90’s [8], extended with a gradient plus error analysis [9] and nonsmooth version [43]. An overview is given in [10], see also [11]. Most convergence analyses are qualitative and limited to convex objectives, only few rates are available. The prescribed step size strategy in Algorithm 1 is directly inspired from these works. We analyse the incremental method as a perturbed gradient method, a view which was exploited in [9, 38] and in distributed settings, see for example [39, 40, 34, 47].

The idea of the adaptive step size is taken from the Adagrad algorithm introduced in [26]. Analysis of such algorithms for nonconvex objectives was proposed in [36, 55, 3, 25] in the stochastic and smooth setting. To our knowledge the combination of adaptive step sizes with incremental methods has not been considered. We use the “scalar step variant” of the Adagrad, called Adagrad-norm in [55] or global step size in [36], in contrast with the originally proposed coordinatewise step sizes analysed in [25]. The original Adagrad algorithm has a power $1/2$ in the denominator which we replaced by $1/3$ in order to obtain faster rates, taking advantage of smoothness and the finite sum structure.

It has been a longstanding open question in machine learning to investigate the advantages of random permutations compared to vanilla SGD [19, 50]. The main motivation is that random permutations often outperforms with replacement sampling despite the absence of theory to explain this observation. The topic is still active and recent progresses has been made in the strongly convex setting, see [56, 29, 48, 54] and reference therein. Rather than studying the superiority of random permutations in the nonconvex setting, we consider the more modest goal of proving convergence guaranties for such strategies. This is achieved by following a perturbed iterate view, a strategy which provides worst case guaranties which are of a different nature compared to average case or almost sure guaranties commonly obtained for stochastic approximation algorithms. The obtained rates have a worse dependency in n compared to SGD but are asymptotically faster than the best known rate for SGD. Similar complexity estimates were obtained in [45, 41] for prescribed step size strategies. To our knowledge, the adaptive variant has not been treated.

Our nonsmooth convergence analysis relies on the ODE method, see [37] with many subsequent developments [6, 33, 7, 17, 3]. In particular we build upon a nonsmooth ODE formulation, differential inclusions [22, 2]. This was used in [23] to analyse the stochastic subgradient algorithm in nonconvex settings using the subgradient projection formula [12]. In the nonsmooth world the backpropagation algorithm [53] used in deep learning suffers from inconsistent behaviors and may not provide subgradient of any kind [31, 32]. We use the recently introduced tool of conservative fields and path differentiable function [14] capturing the full complexity of backpropagation oracles. Our proof essentially relies on the notion of Asymptotic Pseudo Trajectory (APT) for differential inclusions [5, 7].

1.5 Preliminary results

It is important to emphasize that in Algorithm 1, the adaptive step strategy is a special case of the prescribed step strategy. Hence our analysis will start by general considerations for the prescribed step strategy followed by specific considerations to the adaptive steps. We start with a simple claim whose proof is given in appendix A and provides a bound on the length of the steps taken by the algorithm.

Claim 1 *For all $K \in \mathbb{N}$ and all $i = 1, \dots, n$, we have*

$$\max \left\{ \|z_{K,i} - x_K\|^2, \|x_{K+1} - x_K\|^2, \|\hat{z}_{K,i-1} - x_K\|^2 \right\} \leq n \sum_{i=1}^n \alpha_{K,i}^2 \|d(\hat{z}_{K,i-1})\|^2, \quad (2)$$

Throughout this paper, we will work under decreasing step size condition whose meaning is described in the following assumption. We remark that both step size strategies provided in Algorithm 1 comply with this constraint.

Assumption 1 *The sequence $(\alpha_{K,i})_{K \in \mathbb{N}, i \in \{1, \dots, n\}}$ is non increasing with respect to the lexicographic order. That is for all $K \in \mathbb{N}$, $i = 2, \dots, n$, $\alpha_{K,i-1} \geq \alpha_{K,i} \geq \alpha_{K+1,1}$.*

2 Quantitative analysis in the smooth setting

In this section we consider that each f_i has Lipschitz gradient, in which case, d_i is set to be ∇f_i . Note that in this setting, $\nabla F = \frac{1}{n} \sum_{i=1}^n \nabla f_i$ and F also has L -Lipschitz gradient.

Assumption 2 *For $i = 1, \dots, n$,*

- $f_i: \mathbb{R}^p \rightarrow \mathbb{R}$ is an M_i Lipschitz functions and $d_i: \mathbb{R}^p \mapsto \mathbb{R}^p$ is such that for all $x \in \mathbb{R}^p$, $\|d_i(x)\| \leq M_i$. We let $M = \sqrt{\frac{1}{n} \sum_{i=1}^n M_i^2}$. Note that in this case F is M -Lipschitz using Lemma 2 in appendix C which allows to bound $\|\sum_{i=1}^n d_i\|$.
- f_i is continuously differentiable with L_i Lipschitz gradient and we set $d_i = \nabla f_i$. We set $L = \frac{1}{n} \sum_{i=1}^n L_i$. Note that in this case F has L -Lipschitz gradient as shown in Claim 6.

The technical bulk of our analysis is given by the following claim whose proof is provided in appendix A. Note that this result holds deterministically and independantly of the considered step size strategy.

Claim 2 *Under Assumptions 1 and 2, for all $K \in \mathbb{N}$, $K \geq 1$, setting $\alpha_K = \alpha_{K-1,n}$ and $\alpha_0 = \delta^{-1/3} \geq \alpha_{0,1}$, denoting by $(\cdot)_+$ the positive part, we have*

$$\begin{aligned} & F(x_{K+1}) - F(x_K) + \frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 \\ & \leq \left(\alpha_K L^2 n^2 + \left(\frac{Ln}{2} - \frac{1}{2\alpha_K} \right)_+ \right) \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2 + \alpha_K M^2 \sum_{i=1}^n \left(1 - \frac{\alpha_{K,i}^3}{\alpha_K^3} \right). \end{aligned} \quad (3)$$

2.1 Prescribed step size without knowledge of L

The following holds under Assumption for Algorithm 1.

Corollary 1 *If the step size is constant, $\alpha_{K,i} = \alpha/n$ for all $K \in \mathbb{N}$, $i = 1 \dots, n$, we have*

$$\min_{K=0,\dots,N} \|\nabla F(x_K)\|^2 \leq \frac{2(F(x_0) - F^*)}{(N+1)\alpha} + 2 \left(\alpha L^2 M^2 + \frac{LM^2}{2} \right) \alpha$$

Corollary 2 *If the step size is decreasing $\alpha_{K,i} = 1/(n\sqrt{K+1})$, for all $K \in \mathbb{N}$, $i = 1 \dots, n$, then*

$$\min_{K=1,\dots,N} \|\nabla F(x_K)\|^2 \leq \frac{1}{\sqrt{N+1}-1} \left(F(x_0) - F^* + \left(L^2 M^2 + \frac{LM^2}{2} \right) (1 + \log(N+1)) \right)$$

2.2 Prescribed step sizes based on L

The following hold under Assumption 2 for Algorithm 1.

Corollary 3 *If the step size is constant, $\alpha_{K,i} = \alpha/n$, with $\alpha \leq 1/L$, then for all $K \in \mathbb{N}$, $i = 1 \dots, n$, we have*

$$\min_{K=0,\dots,N} \|\nabla F(x_K)\|^2 \leq \frac{2(F(x_0) - F^*)}{(N+1)\alpha} + 2\alpha^2 L^2 M^2$$

Corollary 4 *If the step size is decreasing $\alpha_{K,i} = 1/(Ln(K+1)^{1/3})$, for all $K \in \mathbb{N}$, $i = 1 \dots, n$, then*

$$\min_{K=1,\dots,N} \|\nabla F(x_K)\|^2 \leq \frac{2}{3((N+1)^{2/3}-1)} \left(L(F(x_0) - F^*) + M^2 (1 + \log(N+1)) \right)$$

2.3 Adaptive step size

The following hold under Assumption 2 for Algorithm 1.

Corollary 5 *If we consider the adaptive step size strategy with $\beta = n^2$ and $\delta = n^3$, then*

$$\begin{aligned} & \min_{K=0,\dots,N} \|\nabla F(x_K)\|^2 \\ \leq & 2(M^2 + 1)^{1/3} \frac{F(x_0) - F^* + \left(L^5 + \frac{L^4}{2} \right) + \left(\frac{L^2}{2} (1 + M)^{1/3} + M^2 \right) \log(1 + M^2(N+1))}{(N+1)^{2/3}}. \end{aligned}$$

2.4 Discussion on the obtained convergence rates

All the complexity estimates described in Section 2.2 are given in terms of K , which is the number of epochs. In particular, there is no dependency in the size of the sum n or in the dimension p beyond problem constants L and M . The work presented in [30], see also Theorem 1 in [49], ensures that the convergence rate of “with replacement” SGD applied to problem 1 is of order $O(1/\sqrt{k})$ under the same assumptions as ours, where k is the number of stochastic iteration (typically n times bigger than the number of epochs). From this perspective, the dependency in n is unfavorable as our rates are in terms of number of epochs rather than number of iterations which is customary in stochastic settings [18, 42, 20]. One element of explanation is the nature of our perturbed analysis, which is worst case and blind to the order in which elements are chosen, in contrast with average case stochastic analysis usually performed when considering “with replacement” strategies. This is an important issue, since in practice, for example for deep learning problems, only a few epochs are performed on large datasets. Furthermore, SGD naturally accommodates stochastic data augmentation commonly used in deep learning contexts, a property which is not shared by incremental algorithms.

On the other hand, for prescribed step size and adaptive step size, the convergence rate is of the order of $K^{-2/3}$ which is asymptotically faster than SGD which would be of the order $(nK)^{-1/2}$. This result is only based on comparison of upper bounds and holds only asymptotically since the proposed rate gets better for $K \geq n^3$ which is a regime not considered in practical applications. It constitutes an advantage of the proposed incremental scheme but not a proof of its superiority compared to SGD. Similar rates were obtained in [45] and in [41], with an improved dependency in n and weaker boundedness assumptions. In both cases, it is required to know the Lipschitz constant L which is hardly accessible in practice. The adaptive variant removes this requirement while maintaining a similar rate, showing the advantage of adaptive step sizes in this context. Furthermore, the proposed numerical scheme is more versatile than algorithms in [45, 41] as it allows for a unified treatment of certain form of delays such as minibatching or limited asynchronicity. Interestingly, a cube root variation of the adaptive step size was introduced in [24] in combination with coordinatewise updates and momentum. The purpose is however different, it allows to maintain an effective step size in presence of advanced weighting schemes and is analyzed in a convex setting, while the analysis proposed here takes advantage of larger steps (compared to the vanilla Adagrad algorithm) to obtain faster rates in a nonconvex setting.

Regarding our assumptions, Lipschitzity and boundedness of gradients in Assumption 2 are common in the analysis of stochastic gradient schemes in a nonconvex context, see for example [49, Theorem 1] for SGD and more recently for adaptive variants [25] and incremental variants [45, 41]. Note that in the stochastic approximation context, proxies are often used for these assumptions, requiring them only to hold on the whole sum in (1) rather than on each element. This is often complemented by a uniformly bounded variance assumption, see for example [55]. In finite sum contexts, all these assumptions are very close in nature, as smoothness of the sum directly relates to smoothness of its components. It is worth mentioning that these boundedness assumptions could be relaxed to hold only locally if it is assumed that the sequence remains bounded.

2.5 Proofs for the obtained complexity estimates

Proof of Corollary 1: The considered step size complies with Assumption 2 so that Claim 2 applies. Fix $K \in \mathbb{N}$, fix $\alpha_{K,i} = \alpha_K$ for all $i = 1, \dots, n$, we have $1 - \frac{\alpha_{K,i}^3}{\alpha_K^3} = 0$. Combining with Claim 2, using $\alpha_K \leq \alpha_0$, for all $K \in \mathbb{N}$, we have

$$\frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 \leq F(x_K) - F(x_{K+1}) + \left(\alpha_0 L^2 M^2 n + \frac{LM^2}{2} \right) n^2 \alpha_K^2.$$

Summing for $K = 0, \dots, N$ and dividing by $\sum_{K=0}^N n\alpha_K$, we obtain

$$\min_{K=0, \dots, N} \|\nabla F(x_K)\|^2 \leq \frac{2}{\sum_{K=0}^N n\alpha_K} \left(F(x_0) - F^* + \left(\alpha_0 L^2 M^2 n + \frac{LM^2}{2} \right) \sum_{K=0}^N n^2 \alpha_K^2 \right) \quad (4)$$

Choosing constant step α/n for $\alpha > 0$, we obtain

$$\min_{K=0, \dots, N} \|\nabla F(x_K)\|^2 \leq \frac{2(F(x_0) - F^*)}{(N+1)\alpha} + 2 \left(\alpha L^2 M^2 + \frac{LM^2}{2} \right) \alpha$$

□

Proof of Corollary 2: The considered step size complies with Assumption 2 so that Claim 2 applies. In this setting (4) is still valid. Choosing $\alpha_K = \frac{1}{n\sqrt{K+1}}$, we have

$$\begin{aligned} \sum_{K=0}^N n\alpha_K &\geq \int_{t=0}^{t=N+1} \frac{1}{\sqrt{t+1}} dt \geq 2 \left(\sqrt{N+1} - 1 \right) \\ \sum_{K=0}^N n^2 \alpha_K^2 &\leq \left(1 + \sum_{K=1}^N \frac{1}{K+1} \right) \leq \left(1 + \int_{t=0}^{t=N} \frac{dt}{t+1} \right) = \frac{1}{n^2} (1 + \log(N+1)) \end{aligned}$$

and we obtain in (4)

$$\min_{K=1, \dots, N} \|\nabla F(x_K)\|^2 \leq \frac{1}{\sqrt{N+1} - 1} \left(F(x_0) - F^* + \left(L^2 M^2 + \frac{LM^2}{2} \right) (1 + \log(N+1)) \right)$$

□

Proof of Corollary 3: The considered step size complies with Assumption 2 so that Claim 2 applies. Fix $K \in \mathbb{N}$, fix $\alpha_{K,i} = \alpha_K$ for all $i = 1, \dots, n$, we have $1 - \frac{\alpha_{K,i}^3}{\alpha_K^3} = 0$. Combining with Claim 2, we have, using the fact that $\alpha_K \leq 1/(Ln)$,

$$\frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 \leq F(x_K) - F(x_{K+1}) + L^2 M^2 n^3 \alpha_K^3.$$

Summing for $K = 0, \dots, N$ and dividing by $\sum_{K=0}^N n\alpha_K$, we obtain

$$\min_{K=0, \dots, N} \|\nabla F(x_K)\|^2 \leq \frac{2}{\sum_{K=0}^N n\alpha_K} \left(F(x_0) - F^* + L^2 M^2 \sum_{K=0}^N n^3 \alpha_K^3 \right) \quad (5)$$

Choosing constant step α/n for $\alpha > 0$, we obtain

$$\min_{K=0,\dots,N} \|\nabla F(x_K)\|^2 \leq \frac{2(F(x_0) - F^*)}{(N+1)\alpha} + 2\alpha^2 L^2 M^2$$

□

Proof of Corollary 4: The considered step size complies with Assumption 2 so that Claim 2 applies. In this setting (5) is still valid. Indeed, choosing $\alpha_K = \frac{1}{Ln(K+1)^{1/3}}$, we have $\alpha_K \leq 1/(Ln)$ for all $K \in \mathbb{N}$. Furthermore,

$$\begin{aligned} \sum_{K=0}^N n\alpha_K &\geq \int_{t=0}^{t=N+1} \frac{1}{L(t+1)^{1/3}} dt \geq \frac{3}{2L} ((N+1)^{2/3} - 1) \\ \sum_{K=0}^N n^3 \alpha_K^3 &\leq \frac{1}{L^3} \left(1 + \sum_{K=1}^N \frac{1}{K+1} \right) \leq \frac{1}{L^3} \left(1 + \int_{t=0}^{t=N} \frac{dt}{t+1} \right) = \frac{1}{L^3} (1 + \log(N+1)) \end{aligned}$$

and we obtain in (5)

$$\min_{K=1,\dots,N} \|\nabla F(x_K)\|^2 \leq \frac{2}{3((N+1)^{2/3} - 1)} (L(F(x_0) - F^*) + M^2 (1 + \log(N+1)))$$

□

Proof of Corollary 5: The considered step size complies with Assumption 2 so that Claim 2 applies. We write for all $K \in \mathbb{N}$ and all $i = 1, \dots, n$, $\alpha_K = v_K^{-1/3}$. Let us start with the following.

Claim 3 For all $K \in \mathbb{N}$ and all $i = 1, \dots, n$

$$1 - \frac{\alpha_{K,i}^3}{\alpha_K^3} \leq \beta \sum_{j=1}^n \frac{\|d_j(\hat{z}_{K,j-1})\|^2}{v_{K,j}} \quad (6)$$

Proof of claim 3: Fix $K \in \mathbb{N}$ and i in $1, \dots, n$, we have

$$v_K \leq v_{K,i} = v_K + \beta \sum_{j=1}^i \|d_j(\hat{z}_{K,j-1})\|^2.$$

From this we deduce, using the fact that $v_{K,j}$ is non decreasing in j ,

$$1 - \frac{\alpha_{K,i}^3}{\alpha_K^3} = \frac{v_{K,i} - v_K}{v_{K,i}} = \frac{\beta \sum_{j=1}^i \|d_j(\hat{z}_{K,j-1})\|^2}{v_{K,i}} \leq \beta \sum_{j=1}^i \frac{\|d_j(\hat{z}_{K,j-1})\|^2}{v_{K,j}} \leq \beta \sum_{j=1}^n \frac{\|d_j(\hat{z}_{K,j-1})\|^2}{v_{K,j}},$$

□

Combining Claim 2 and Claim 3, we have for all $K \in \mathbb{N}$, using $\delta^{-1/3} \geq \alpha_K$

$$\begin{aligned} \frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 &\leq F(x_K) - F(x_{K+1}) + \left(\alpha_K L^2 n^2 + \left(\frac{Ln}{2} - \frac{1}{2\alpha_K} \right)_+ \right) \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2 \\ &\quad + M^2 n \frac{\beta}{\delta^{1/3}} \sum_{j=1}^n \alpha_{K,j}^3 \|d_j(\hat{z}_{K,j-1})\|^2 \end{aligned} \quad (7)$$

We will consider the following notation $\bar{\alpha} = \frac{1}{Ln}$, we have

$$\left(\alpha_K L^2 n^2 + \left(\frac{Ln}{2} - \frac{1}{2\alpha_K} \right)_+ \right) \leq \alpha_K L^2 n^2$$

if and only if $\alpha_K \leq \bar{\alpha}$. Set \bar{K} , the first index K such that $\alpha_K \leq \bar{\alpha}$. For all $K \leq \bar{K} - 1$, we have $1/\delta^{1/3} \geq \alpha_K = v_K^{-1/3} > \bar{\alpha}$. Fix $N \leq \bar{K} - 1$, summing the second term of (7) for $K = 0 \dots N$, we have

$$\begin{aligned} & \sum_{K=0}^N \left(\alpha_K L^2 n^2 + \left(\frac{Ln}{2} - \frac{1}{2\alpha_K} \right)_+ \right) \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2 \\ & \leq \left(\frac{L^2 n^2}{\delta} + \frac{Ln}{2\delta^{2/3}} \right) \sum_{K=0}^N \sum_{j=1}^n \|d_j(\hat{z}_{K,j-1})\|^2 \\ & \leq \left(\frac{L^2 n^2}{\beta\delta} + \frac{Ln}{2\beta\delta^{2/3}} \right) v_{\bar{K}} \\ & \leq \left(\frac{L^2 n^2}{\beta\delta} + \frac{Ln}{2\beta\delta^{2/3}} \right) \frac{1}{\bar{\alpha}^3} \end{aligned} \quad (8)$$

Now, choosing $N \geq \bar{K}$, summing the same quantity for $K \geq \bar{K}$, we have using the definition of $\bar{\alpha}$

$$\begin{aligned} & \sum_{K=\bar{K}}^N \left(\alpha_K L^2 n^2 + \left(\frac{Ln}{2} - \frac{1}{2\alpha_K} \right)_+ \right) \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2 \\ & \leq \sum_{K=\bar{K}}^N \alpha_K L^2 n^2 \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2 \\ & \leq L^2 n^2 \sum_{j=1}^n \sum_{K=0}^N \frac{\alpha_K}{\alpha_{K,j}} \alpha_{K,j}^3 \|d_j(\hat{z}_{K,j-1})\|^2 \\ & \leq L^2 n^2 (1 + \beta n M / \delta)^{1/3} \sum_{j=1}^n \sum_{K=0}^N \alpha_{K,j}^3 \|d_j(\hat{z}_{K,j-1})\|^2 \end{aligned} \quad (9)$$

where the last identity follows because for all $K \in \mathbb{N}$ and $j = 1 \dots n$,

$$\frac{\alpha_K^3}{\alpha_{K,j}^3} = \frac{v_{K,j}}{v_K} = \frac{v_K + \beta \sum_{i=1}^j \|d_i(\hat{z}_{K,i-1})\|^2}{v_K} \leq 1 + \frac{\beta n M}{\delta}$$

Combining (8) and (9), for any $N \in \mathbb{N}$, independently of its position relative to \bar{K} (and even if $\bar{K} = +\infty$), we have

$$\begin{aligned} & \sum_{K=0}^N \left(\alpha_K L^2 n^2 + \left(\frac{Ln}{2} - \frac{1}{2\alpha_K} \right)_+ \right) \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2 \\ & \leq \left(\frac{L^2 n^2}{\beta\delta} + \frac{Ln}{2\beta\delta^{2/3}} \right) \frac{1}{\bar{\alpha}^3} + L^2 n^2 (1 + \beta n M / \delta)^{1/3} \sum_{j=1}^n \sum_{K=0}^N \alpha_{K,j}^3 \|d_j(\hat{z}_{K,j-1})\|^2 \end{aligned} \quad (10)$$

Given $N \in \mathbb{N}$, we may sum (7) for $K = 0 \dots, N$ combined with (10) to obtain

$$\begin{aligned} \sum_{K=0}^N \frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 &\leq F(x_0) - F(x_N) + \left(\frac{L^2 n^2}{\beta\delta} + \frac{Ln}{2\beta\delta^{2/3}} \right) \frac{1}{\bar{\alpha}^3} \\ &\quad + \left(L^2 n^2 (1 + \beta n M / \delta)^{1/3} + M^2 n \frac{\beta}{\delta^{1/3}} \right) \sum_{j=1}^n \sum_{K=0}^N \alpha_{K,j}^3 \|d_j(\hat{z}_{K,j-1})\|^2 \end{aligned} \quad (11)$$

Now, we use the lexicographic order on pairs of integers, $(a, b) \leq (c, d)$ if $a < c$ or $a = c$ and $b \leq d$. From Lemma 3 in appendix C, we have

$$\begin{aligned} \sum_{K=0}^N \sum_{i=1}^n \alpha_{K,i}^3 \|d_i(\hat{z}_{K,i-1})\|^2 &= \sum_{(K,i) \leq (N,n)} \frac{\|d_i(\hat{z}_{K,i-1})\|^2}{\delta + \beta \sum_{(k,j) \leq (K,i)} \|d_j(\hat{z}_{k,j-1})\|^2} \\ &\leq \frac{1}{\beta} \log \left(1 + \frac{\beta \sum_{(K,i) \leq (N,n)} \|d_i(\hat{z}_{K,i-1})\|^2}{\delta} \right) \leq \frac{1}{\beta} \log \left(1 + \frac{\beta n M^2 (N+1)}{\delta} \right), \end{aligned} \quad (12)$$

where the first inequality follows by applying Lemma 3, noticing that we sum over $(N+1)n$ instances and that $\sum_{i=1}^n \|d_i\|^2 \leq nM^2$. We remark that for all $K \in \mathbb{N}$, $\alpha_K \geq (Kn\beta M^2 + \delta)^{-1/3}$. Combining (11) and (12), we obtain

$$\begin{aligned} &\frac{n(N+1)(Nn\beta M^2 + \delta)^{-1/3}}{2} \min_{K=0, \dots, N} \|\nabla F(x_K)\|^2 \\ &\leq F(x_0) - F^* + \left(\frac{L^2 n^2}{\beta\delta} + \frac{Ln}{2\beta\delta^{2/3}} \right) \frac{1}{\bar{\alpha}^3} + \left(\frac{L^2 n^2}{\beta} (1 + \beta n M / \delta)^{1/3} + \frac{M^2 n}{\delta^{1/3}} \right) \log \left(1 + \frac{\beta n M^2 (N+1)}{\delta} \right). \end{aligned} \quad (13)$$

Combining (13) with $\bar{\alpha} = 1/(Ln)$, and choosing $\beta = n^2$ and $\delta = n^3$, we obtain

$$\begin{aligned} &\frac{(N+1)}{2(NM^2 + 1)^{1/3}} \min_{K=0, \dots, N} \|\nabla F(x_K)\|^2 \\ &\leq F(x_0) - F^* + \left(L^5 + \frac{L^4}{2} \right) + \left(\frac{L^2}{2} (1 + M)^{1/3} + M^2 \right) \log(1 + M^2(N+1)). \end{aligned}$$

The rest follows by noticing that $\frac{(N+1)}{2(NM^2+1)^{1/3}} \geq \frac{(N+1)}{2((N+1)(M^2+1))^{1/3}} = \frac{(N+1)^{2/3}}{2(M^2+1)^{1/3}}$ \square

3 Qualitative analysis for nonsmooth objectives

In this section we consider nonsmooth objectives such as typical losses arising when training deep networks. Our analysis will be performed under the following standing assumption.

Assumption 3 *In addition to Assumption 1, assume that*

$$\sum_{K=0}^{\infty} \alpha_{K,1} = +\infty, \quad \text{and} \quad \alpha_{K,1} \xrightarrow{K \rightarrow \infty} 0, \quad \text{and} \quad \frac{\alpha_{K,1}}{\alpha_{K,n}} \xrightarrow{K \rightarrow \infty} 1. \quad (14)$$

We follow the ODE approach, our arguments closely follow those developed in [7]. We start by defining a continuous time piecewise affine interpolant of the sequence.

Definition 1 For all $K \in \mathbb{N}$, we let $\tau_K = \sum_{k=0}^K \sum_{i=1}^n \alpha_{k,i}$. We fix the sequence given by Algorithm (1) and consider the associated Lipschitz interpolant such that $\mathbf{w}: \mathbb{R}^+ \mapsto \mathbb{R}^p$, such that $\mathbf{w}(\tau_K) = x_K$ for all $K \in \mathbb{N}$ and the interpolation is affine on (τ_K, τ_{K+1}) for all $K \in \mathbb{N}$.

3.1 Differential inclusion setting

The main argument in this Section is connecting the continuous time interpolant in Definition 1 and continuous dynamics. The continuous time counterpart of Algorithm 1, is $\dot{\mathbf{x}} = \frac{-1}{n} \sum_{i=1}^n d_i(\mathbf{x})$, for which the right hand side is not continuous, classical Cauchy-Lipschitz type theorems for existence of solutions cannot be applied. We need to resort to a continuous extension of the right hand side, which becomes set valued, providing a weaker notion of solution. We use the recently introduced notion of conservativity [14] which captures the complexity of automatic differentiation oracles in nonsmooth settings [15]. Recall that the set valued map D is conservative for the locally Lipschitz function f , if it has a closed graph and for any locally Lipschitz curve $\mathbf{x}: [0, 1] \mapsto \mathbb{R}^p$ and almost all $t \in [0, 1]$

$$\frac{d}{dt}f(\mathbf{x}(t)) = \langle v, \dot{\mathbf{x}}(t) \rangle, \quad \forall v \in D(\mathbf{x}(t)). \quad (15)$$

This is the counterpart to $\frac{d}{dt}f(\mathbf{x}(t)) = \langle \nabla f(x(t)), \dot{\mathbf{x}}(t) \rangle$ for any C^1 function f and any C^1 curve \mathbf{x} . This property is known as the chain rule of subdifferential inclusions, see for example [23]. The main specificity is that the property holds for almost all t due to the fact that we have nondifferentiable objects, and for all possible choices in D which is set valued, again due to nondifferentiability. As shown in [14], this ensures that for any such curve, one has

$$f(\mathbf{x}(1)) - f(\mathbf{x}(0)) = \int \max_{v \in D(\mathbf{x}(t))} \langle v, \dot{\mathbf{x}}(t) \rangle dt = \int \min_{v \in D(\mathbf{x}(t))} \langle v, \dot{\mathbf{x}}(t) \rangle dt,$$

where the integral is understood in the Lebesgue sense.

Assumption 4 For $i = 1, \dots, n$, we let D_i be a conservative field for f_i with $\max_{v \in D_i(x)} \|v\| \leq M$ for all $x \in \mathbb{R}^p$ and $d_i: \mathbb{R}^p \mapsto \mathbb{R}^p$ is measurable such that for all $x \in \mathbb{R}^p$, $d_i(x) \in D_i(x)$. We set $D = \text{conv} \left(\frac{1}{n} \sum_{i=1}^n D_i \right)$. Since conservativity is preserved under addition [14, Corollary 4] D is conservative for F , furthermore it has convex compact values and a closed graph. We set crit_F to be the set of $x \in \mathbb{R}^p$ such that $0 \in D(x)$.

Main examples in deep learning: If each f_i , $i = 1, \dots, n$ is the loss associated to a sample point and a neural network architecture, assuming that f_i is defined using a compositional formula involving piecewise polynomials, logarithms and exponentials (which covers most of deep network architectures), then the Clarke subgradient [22] is a

conservative field for f_i . Recall that the Clarke subgradient extends the notion of convex subgradient to nonconvex locally Lipschitz functions. This was proved in [23] using the projection formula in [14, 15], see also [21, 14]. In deep learning context, backpropagation may fail to provide Clarke subgradients in nonsmooth contexts [31, 32]. Nonetheless, it was shown in [14] that backpropagation computes a conservative field. Hence our analysis applies to training of deep networks using a backpropagation oracle such as the ones implemented in [1, 46].

Definition 2 *A solution to the differential inclusion*

$$\dot{\mathbf{x}} \in -D(\mathbf{x})$$

with initial point $x \in \mathbb{R}^p$ is a locally Lipschitz mapping $\mathbf{x}: \mathbb{R} \mapsto \mathbb{R}^p$ such that $\mathbf{x}(0) = x$ and for almost all $t \in \mathbb{R}$, $\dot{\mathbf{x}}(t) \in -D(\mathbf{x}(t))$. We denote by S_x the set of such solutions and by S the set of all solutions with any initialization.

Standard results in this field [2, Chapter 2, Theorem 3] ensure that, since D has closed graph and compact convex values, for any $x \in \mathbb{R}^p$ the set S_x is nonempty, note that it could be non unique.

3.2 Main result

The following notion was introduced in [7], see also [5]. It captures the fact that a continuous trajectory is a solution to the differential inclusion in Definition 2 asymptotically. Note that we let the initialization free in the next definition, this is necessary to apply [7, Theorem 4.1].

Definition 3 (Asymptotic pseudo trajectory) *A continuous function $\mathbf{z}: \mathbb{R}_+ \mapsto \mathbb{R}^p$ is an asymptotic pseudotrajectory (APT), if for all $T > 0$,*

$$\liminf_{t \rightarrow \infty} \sup_{\mathbf{x} \in S} \sup_{0 \leq s \leq T} \|\mathbf{z}(t+s) - \mathbf{x}(s)\| = 0.$$

Claim 4 *Under Assumptions 1, 3 and 4, assume that $(x_K)_{K \in \mathbb{N}}$ produced by Algorithm 1 with prescribed step size is bounded. Then the interpolant \mathbf{w} given in Definition 1 is an asymptotic pseudo trajectory as described in Definition 3.*

The proof relies on Lemma 1 which shows that the iterates produced by the algorithm satisfy a perturbed differential inclusion. The technical bulk of the proof is in Theorem 1 which shows that perturbed differential inclusions are asymptotic pseudo trajectories. These results are described in Section 3.3, the presentation and main arguments follow the ideas presented in [7]. In order to deduce convergence of Algorithm 1 from the Asymptotic pseudo trajectory property, we need the following Morse-Sard assumption. We stress that for deep network involving piecewise polynomials, logarithms and exponentials, this assumption is satisfied for both the Clarke subgradient and the backpropagation oracle [12, 23, 14].

Assumption 5 *The function F and D are such that $F(\text{crit}_F)$, does not contain any open interval, where crit_F is given in Assumption 4 and contains all $x \in \mathbb{R}^p$, with $0 \in D(x)$.*

Corollary 6 *Under Assumptions 1, 3 and 4, assume that $(x_K)_{K \in \mathbb{N}}$ produced by Algorithm 1 with prescribed step size is bounded and that Assumption 5 holds. Then $F(x_K)$ converges to a critical value of F as $K \rightarrow \infty$ and all accumulation points of the sequence are critical points for D .*

Proof : Let $\mathbf{x}: \mathbb{R}^p \mapsto \mathbb{R}$ be a solution to the differential inclusion described in Definition 2. Then using conservativity in (15), for almost all $t \in \mathbb{R}_+$, we have

$$\frac{d}{dt}F(\mathbf{x}(t)) = - \min_{v \in D(\mathbf{x}(t))} \|v\|^2$$

Hence F is a Lyapunov function for the system: it decreases along trajectory, strictly outside crit_F . Using Claim 4, \mathbf{w} is an APT. Combining Assumption 5 with Proposition 3.27 and Theorem 4.3 in [7], all limit points of \mathbf{w} are contained in crit_F and F is constant on this set, that is $F(\mathbf{w}(t))$ converges as $t \rightarrow \infty$. \square

Corollary 7 *Under Assumption 4, assume that $(x_K)_{K \in \mathbb{N}}$ produced by Algorithm 1 with adaptive step size is bounded and that Assumption 5 holds. Then $F(x_K)$ converges to a critical value of F as $K \rightarrow \infty$ and all accumulation points of the sequence are critical points for D .*

Proof : If v_K converges, this means that all d_i go to 0, and all partial increments also vanish asymptotically due to Claim 1. Call the set of accumulation points $\Omega \subset \mathbb{R}^p$. Ω forms a compact connected subset of crit_F , see [13, Lemma 3.5, (iii)] for details. By continuity of F , the $F(\Omega)$ is a connected subset of \mathbb{R} , that is an interval. By Morse-Sard assumption 5 it is a singleton which proves the claim. Assume otherwise that v_K diverges to $+\infty$ as $K \rightarrow \infty$, in this case, the step size goes to 0. We have

$$v_K \leq v_{K+1} \leq v_K + nM$$

which shows that $v_{K+1}/v_K \rightarrow 1$ as $K \rightarrow \infty$, and $\sum_{K \in \mathbb{N}} \alpha_{K,1} = +\infty$ so that Assumptions 1 and 3 are valid and Corollary 7 applies. \square

3.3 Proof of the main result

We extend and adapt the arguments of [7].

Definition 4 (Local extension) *For any $\gamma > 0$, and any $x \in \mathbb{R}^p$, we let D^γ be the following local extension of D*

$$D^\gamma(x) = \left\{ y \in \mathbb{R}^p, y \in \frac{1}{n} \sum_{i=1}^n \lambda_i D_i(x_i), \|x - x_i\| \leq \gamma, |\lambda_i - 1| \leq \gamma, i = 1, \dots, n \right\}.$$

Note that $\lim_{\gamma \rightarrow 0} D^\gamma(x) = \frac{1}{n} \sum_{i=1}^n D_i(x)$ by graph closedness of each D_i in Assumption 4.

Definition 5 (Perturbed differential inclusion) *A locally Lipschitz path $\mathbf{x}: \mathbb{R}_+ \mapsto \mathbb{R}^p$ satisfies the perturbed differential inclusion if there exists a function $\gamma: \mathbb{R}_+ \mapsto \mathbb{R}_+$ with $\lim_{t \rightarrow \infty} \gamma(t) = 0$, such that for almost all $t \geq 0$*

$$\dot{\mathbf{x}}(t) \in -D^{\gamma(t)}(\mathbf{x}(t))$$

Lemma 1 *The interpolated trajectory \mathbf{w} given in Definition 1 satisfies the perturbed differential inclusion in Definition 5.*

Proof : The interpolated trajectory is piecewise affine so it is locally Lipschitz and differentiable almost everywhere. For each $K \in \mathbb{N}$ and $i = 1, \dots, n$, we have using Claim 1

$$\|x_K - \hat{z}_{K,i-1}\| \leq n\alpha_{K,1}M. \quad (16)$$

Furthermore, for all $t \in (\tau_K, \tau_{K+1})$,

$$\dot{\mathbf{w}}(t) = - \sum_{i=1}^n \alpha_{K,i} d_i(\hat{z}_{K,i-1}) / (\tau_{K+1} - \tau_K) = - \frac{1}{n} \sum_{i=1}^n \lambda_i d_i(\hat{z}_{K,i-1}), \quad (17)$$

where for all $i = 1, \dots, n$, using $\alpha_{K,i} \leq \alpha_{K,1}$ and $\tau_{K+1} - \tau_K = \sum_{i=1}^n \alpha_{K,i} \geq n\alpha_{K,n}$,

$$\lambda_i = \frac{n\alpha_{K,i}}{\tau_{K+1} - \tau_K} \leq n \frac{\alpha_{K,1}}{n\alpha_{K,n}} = \frac{\alpha_{K,1}}{\alpha_{K,n}}. \quad (18)$$

Hence combining (16) and (17), we may consider $\gamma(t) = \max \left\{ n\alpha_{K,1}M, \left| 1 - \frac{\alpha_{K,1}}{\alpha_{K,n}} \right| \right\}$ for all $t \in (\tau_K, \tau_{K+1})$ which satisfies the desired hypothesis. \square

The following result is the main technical part of this section. The proof follows that of [7, Theorem 4.2] and is provided in Appendix B.

Theorem 1 *Let \mathbf{z} be a perturbed differential inclusion trajectory as given in Definition 5. Then \mathbf{z} is an asymptotic pseudotrajectory as described in Definition 3.*

3.4 Discussion of the obtained result

Definition 5 extends the notion of approximate differential inclusion introduced in [7] to the finite sum setting. Indeed, the definition proposed in [7] coincides with ours when $n = 1$. We add the flexibility to choose different approximation points for each elements of the sum which, in turn, allows to conclude regarding the output of the algorithm. A more general result was described in [13] in a more abstract form. It is interesting to notice that the differential inclusion approach was developed to analyze stochastic approximation algorithms because of the difficulty caused by the addition of random noise. The proposed analysis suggests that this approach is also useful to analyse deterministic algorithms as ours.

The obtained convergence result is qualitative and completely mirrors what is obtained for SGD under similar assumptions at this level of generality [23, 14]. In terms of assumptions, analysis of stochastic approximation requires that the step size decay is proportioned to concentration of the noise. For example, under uniformly bounded variance, step sizes should be square summable. Such assumptions ensure that perturbations of dynamics induced by the noise are summable, and therefore negligible in the limit, see for example [6] for a discussion. Such an assumption is not required by our deterministic approach, the only required assumption is that the step size goes to zero in the limit and that the steps remain of the same order within an epoch.

All the obtained results hold under the assumption that the trajectory remains bounded. This is a strong assumption which is difficult to check a priori given a problem of the form (1). This assumption is common in the analysis of stochastic approximation algorithms [7, 23] and we are not aware of easy sufficient condition which ensures that this is the case. A simple work around would be to add a projection step on a compact convex set at the end of each epoch. This would correspond to a constrained optimization problem in place of (1). The considered notion of approximate differential inclusion in Definition 5 is general enough to include this additional algorithmic step in the analysis, maintaining the qualitative convergence result without requiring the boundedness assumption, which would be automatically fulfilled.

4 Conclusion

We have introduced a flexible algorithmic framework for finite sums and proposed convergence guaranties in smooth and nonsmooth settings under assumptions which are qualitatively similar as in the litterature on stochastic gradient descent for such problems. The obtained result rely on a perturbed iterate analysis and are valid in a worse case sense, they have therefore a quite different nature compared to guaranties obtained for stochastic approximation algorithms. In the smooth setting we obtain quantitative rates which have worse dependency in n but are asymptotically faster. The resulting complexity estimate improves over SGD in the asymptotic regime, but remains weaker for first epochs, a situation which is not uncommon in the analysis of incremental methods [28].

A natural extension of this work would consist in providing proof arguments explaining why random permutations, as implemented in practice, often provide superior results compared to “with replacement sampling” in a nonasymptotic sense. This topic has been extensively studied in the strongly convex setting [19, 50, 56, 29, 48, 54] and it is of interest to extend these ideas to the nonconvex and possibly nonsmooth setting [45, 41]. This will be the subject of future research. Finally, another topic of interest would be to devise variants of the proposed algorithmic scheme with faster convergence rates.

Acknowledgements. The authors acknowledge the support of ANR-3IA Artificial and Natural Intelligence Toulouse Institute, Air Force Office of Scientific Research, Air Force Material Command, USAF, under grant numbers FA9550-19-1-7026, FA9550-18-1-0226, and ANR MasDol. The author would like to thank anonymous referees for their comments which helped improve the relevance of the paper.

References

- [1] Abadi M., Barham P., Chen J., Chen Z., Davis A., Dean J., Devin M., Ghemawat S., Irving G., Isard M., Kudlur M., Levenberg J., Monga R., Moore S., Murray D., Steiner B., Tucker P., Vasudevan V., Warden P., Wicke M., Yu Y. and Zheng X. (2016). Tensorflow: A system for large-scale machine learning. In Symposium on Operating Systems Design and Implementation.
- [2] Aubin, J. P., Cellina, A. (1984). Differential inclusions: set-valued maps and viability theory (Vol. 264). Springer.
- [3] Barakat, A., & Bianchi, P. (2018). Convergence and Dynamical Behavior of the Adam Algorithm for Non Convex Stochastic Optimization. arXiv preprint arXiv:1810.02263.
- [4] Baydin A., Pearlmutter B., Radul A. and Siskind J. (2018). Automatic differentiation in machine learning: a survey. *Journal of machine learning research*, 18(153).
- [5] Benaïm, M., & Hirsch, M. W. (1996). Asymptotic pseudotrajectories and chain recurrent flows, with applications. *Journal of Dynamics and Differential Equations*, 8(1), 141-176.
- [6] Benaïm, M. (1999). Dynamics of stochastic approximation algorithms. In *Séminaire de probabilités XXXIII* (pp. 1-68). Springer, Berlin, Heidelberg.
- [7] Benaïm, M., Hofbauer, J., Sorin, S. (2005). Stochastic approximations and differential inclusions. *SIAM Journal on Control and Optimization*, 44(1), 328-348.
- [8] Bertsekas, D. P. (1997). A new class of incremental gradient methods for least squares problems. *SIAM Journal on Optimization*, 7(4), 913-926.
- [9] Bertsekas, D. P., & Tsitsiklis, J. N. (2000). Gradient convergence in gradient methods with errors. *SIAM Journal on Optimization*, 10(3), 627-642.
- [10] Bertsekas, D. P. (2011). Incremental gradient, subgradient, and proximal methods for convex optimization: A survey. *Optimization for Machine Learning*, 2010(1-38), 3.
- [11] Bertsekas, D. P. (2015). *Convex optimization algorithms*. Belmont: Athena Scientific.
- [12] Bolte, J., Daniilidis, A., Lewis, A., Shiotani, M. (2007). Clarke subgradients of stratifiable functions. *SIAM Journal on Optimization*, 18(2), 556-572.
- [13] Bolte, J., Sabach, S., & Teboulle, M. (2014). Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1-2), 459-494.

- [14] Bolte, J. and Pauwels, E. (2020). Conservative set valued fields, automatic differentiation, stochastic gradient methods and deep learning. *Mathematical Programming*.
- [15] Bolte J., and Pauwels E. (2020). A mathematical model for automatic differentiation in machine learning. In *Conference on Neural Information Processing Systems*.
- [16] Bolte J., Pauwels E. and Rios-Zertuche R. (2020). Long term dynamics of the subgradient method for Lipschitz path differentiable functions. *arXiv preprint arXiv:2006.00098*.
- [17] Borkar, V. (2009). *Stochastic approximation: a dynamical systems viewpoint* (Vol. 48). Springer.
- [18] Bottou L. and Bousquet O. (2008). The tradeoffs of large scale learning. In *Advances in neural information processing systems* (pp. 161-168).
- [19] Bottou L. (2009). Curiously fast convergence of some stochastic gradient descent algorithms. In *Proceedings of the symposium on learning and data science, Paris* (Vol. 8, pp. 2624-2633).
- [20] Bottou L., Curtis F. E. and Nocedal J. (2018). Optimization methods for large-scale machine learning. *Siam Review*, 60(2), 223-311.
- [21] Castera C., Bolte J., Févotte C., Pauwels E. (2019). An Inertial Newton Algorithm for Deep Learning. *arXiv preprint arXiv:1905.12278*.
- [22] Clarke F. H. (1983). *Optimization and nonsmooth analysis*. Siam.
- [23] Davis, D., Drusvyatskiy, D., Kakade, S., Lee, J. D. (2018). Stochastic subgradient method converges on tame functions. *Foundations of Computational Mathematics*.
- [24] Defazio, A., & Jelassi, S. (2021). Adaptivity without Compromise: A Momentumized, Adaptive, Dual Averaged Gradient Method for Stochastic Optimization. *arXiv preprint arXiv:2101.11075*.
- [25] Défossez, A., Bottou, L., Bach, F., & Usunier, N. (2020). On the Convergence of Adam and Adagrad. *arXiv preprint arXiv:2003.02395*.
- [26] Duchi, J., Hazan, E., & Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12, 2121-2159.
- [27] Ekeland, I., & Temam, R. (1976). *Convex analysis and variational problems*. SIAM.
- [28] Gürbüzbalaban Mert, Asuman Ozdaglar and Pablo A. Parrilo. On the convergence rate of incremental aggregated gradient algorithms. *SIAM Journal on Optimization* 27.2 (2017): 1035-1048.

- [29] Gürbüzbalaban, Mert, Asu Ozdaglar, and Pablo A. Parrilo. Why random reshuffling beats stochastic gradient descent. *Mathematical Programming* (2019): 1-36.
- [30] Ghadimi, S. and Lan, G. (2013). Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4), 2341-2368.
- [31] Griewank, A., Walther, A. (2008). Evaluating derivatives: principles and techniques of algorithmic differentiation (Vol. 105). SIAM.
- [32] Kakade, S. M. and Lee, J. D. (2018). Provably correct automatic sub-differentiation for qualified programs. In *Advances in Neural Information Processing Systems* (pp. 7125-7135).
- [33] Kushner H. and Yin, G. G. (2003). Stochastic approximation and recursive algorithms and applications (Vol. 35). Springer Science & Business Media.
- [34] Lan, G., Lee, S., & Zhou, Y. (2018). Communication-efficient algorithms for decentralized and stochastic optimization. *Mathematical Programming*.
- [35] LeCun Y., Bengio Y., Hinton, G. (2015). Deep learning. *Nature*, 521(7553).
- [36] Li, X., & Orabona, F. (2019). On the Convergence of Stochastic Gradient Descent with Adaptive Stepsizes. In *International Conference on Artificial Intelligence and Statistics*.
- [37] Ljung L. (1977). Analysis of recursive stochastic algorithms. *IEEE transactions on automatic control*, 22(4), 551-575.
- [38] Mania, H., Pan, X., Papailiopoulos, D., Recht, B., Ramchandran, K., & Jordan, M. I. (2017). Perturbed iterate analysis for asynchronous stochastic optimization. *SIAM Journal on Optimization*, 27(4), 2202-2229.
- [39] McMahan, B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Artificial Intelligence and Statistics* (pp. 1273-1282).
- [40] Mishchenko, K., Iutzeler, F., Malick, J., & Amini, M. R. (2018). A delay-tolerant proximal-gradient algorithm for distributed learning. In *International Conference on Machine Learning* (pp. 3587-3595).
- [41] Mishchenko, K., Khaled Ragab Bayoumi, A., & Richtárik, P (2020). Random reshuffling: Simple analysis with vast improvements. *Advances in Neural Information Processing Systems*, 33.
- [42] Moulines E. and Bach, F. (2011). Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *Advances in Neural Information Processing Systems* (pp. 451-459).

- [43] Nedic, A., & Bertsekas, D. P. (2001). Incremental subgradient methods for nondifferentiable optimization. *SIAM Journal on Optimization*, 12(1), 109-138.
- [44] Nesterov, Y. (2004). *Introductory lectures on convex optimization: A basic course*. Springer Science & Business Media.
- [45] Nguyen, L. M., Tran-Dinh, Q., Phan, D. T., Nguyen, P. H., & van Dijk, M. (2020). A unified convergence analysis for shuffling-type gradient methods. arXiv preprint arXiv:2002.08246.
- [46] Paszke A., Gross S., Chintala S., Chanan G., Yang E., DeVito Z., Lin Z., Desmaison A., Antiga L. and Lerer A. (2017). Automatic differentiation in pytorch. In NIPS workshops.
- [47] Pu, S., & Nedic A. (2020). Distributed stochastic gradient tracking methods. *Mathematical Programming*, 1-49.
- [48] Rajput, S., Gupta, A. and Papailiopoulos, D. (2020). Closing the convergence gap of SGD without replacement. In *International Conference on Machine Learning* (pp. 7964-7973). PMLR.
- [49] Sashank J. Reddi, Suvrit Sra, Barnabas Poczos and Alexander J. Smola (2016). Fast incremental method for smooth nonconvex optimization. *IEEE Conference on Decision and Control (CDC)*.
- [50] Recht B., and Ré C. (2012). Toward a noncommutative arithmetic-geometric mean inequality: conjectures, case-studies, and consequences. In *Conference on Learning Theory* (pp. 11-1). *JMLR Workshop and Conference Proceedings*.
- [51] Robbins H. and Monro, S. (1951). A stochastic approximation method. *The annals of mathematical statistics*, 400-407.
- [52] Royden, H. L., & Fitzpatrick, P. (1988). *Real analysis*. New York: Macmillan.
- [53] Rumelhart E., Hinton E., Williams J. (1986). Learning representations by back-propagating errors. *Nature* 323:533-536.
- [54] Safran, I. and Shamir, O. (2020). How good is SGD with random shuffling? In *Conference on Learning Theory* (pp. 3250-3284). PMLR.
- [55] Ward, R., Wu, X., & Bottou, L. (2019). Adagrad stepsizes: Sharp convergence over nonconvex landscapes, from any initialization. In *International Conference on Machine Learning*.
- [56] Ying, B., Yuan, K., Vlaski, S. and Sayed, A. H. (2018). Stochastic learning under random reshuffling with constant step-sizes. *IEEE Transactions on Signal Processing*, 67(2), 474-489.
- [57] Zhang, J., Lin, H., Sra, S., & Jadbabaie, A. (2020). On Complexity of Finding Stationary Points of Nonsmooth Nonconvex Functions. arXiv preprint arXiv:2002.04130.

This is the appendix for “Incremental Without Replacement Sampling in Nonconvex Optimization”. We begin with the proof of the first claim of the paper.

Proof of Claim 1: We have for all $K \in \mathbb{N}$ and $i = 1 \dots n$, using the recursion in Algorithm 1,

$$z_{K,i} - x_K = \sum_{j=1}^i \alpha_{K,j} d(\hat{z}_{K,j-1}).$$

Using Lemma 2, we obtain

$$\|z_{K,i} - x_K\|^2 \leq i \sum_{j=1}^i \alpha_{K,i}^2 \|d(\hat{z}_{K,i-1})\|^2 \leq n \sum_{i=1}^n \alpha_{K,i}^2 \|d(\hat{z}_{K,i-1})\|^2.$$

Taking $i = n$, we obtain the second inequality. The result follows for $\hat{z}_{K,i-1}$ because it is in $\text{conv}(z_{K,j})_{j=0}^{i-1}$ and

$$\|\hat{z}_{K,i-1} - x_K\|^2 \leq \max_{z \in \text{conv}(z_{K,j})_{j=0}^{i-1}} \|z - x_K\|^2 = \max_{j=0, \dots, i} \|z_{K,j} - x_K\|^2 \leq n \sum_{i=1}^n \alpha_{K,i}^2 \|d(\hat{z}_{K,i-1})\|^2,$$

where the equality in the middle follows because the maximum of a convex function over a polyhedra is achieved at vertices. \square

A Proofs for the smooth setting

For all $K \in \mathbb{N}$, we let $\alpha_K = \alpha_{K-1,n}$, with $\alpha_0 = \delta^{-1/3} \geq \alpha_{0,1}$.

A.1 Analysis for both step size strategies.

Claim 5 *We have for all $K \in \mathbb{N}$,*

$$\begin{aligned} & \langle \nabla F(x_K), x_{K+1} - x_K \rangle + \frac{1}{2n\alpha_K} \|x_{K+1} - x_K\|^2 \\ & \leq -\frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 + \alpha_K L^2 n^2 \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2 + \alpha_K M^2 \sum_{i=1}^n \left(\frac{\alpha_{K,i}}{\alpha_K} - 1 \right)^2 \end{aligned} \quad (19)$$

Proof of Claim 5: Fix $K \in \mathbb{N}$, we have

$$x_{K+1} - x_K = - \sum_{i=1}^n \alpha_{K,i} d_i(\hat{z}_{K,i-1}) = -\alpha_K \sum_{i=1}^n \frac{\alpha_{K,i}}{\alpha_K} d_i(\hat{z}_{K,i-1}) \quad (20)$$

Recall that $\nabla F(x_K) = \frac{1}{n} \sum_{i=1}^n d_i(x_K)$, combining with (20), we deduce the following

$$\begin{aligned}
& \langle \nabla F(x_K), x_{K+1} - x_K \rangle + \frac{1}{2n\alpha_K} \|x_{K+1} - x_K\|^2 \\
= & \frac{-\alpha_K}{n} \left\langle \sum_{i=1}^n d_i(x_K), \sum_{i=1}^n \frac{\alpha_{K,i}}{\alpha_K} d_i(\hat{z}_{K,i-1}) \right\rangle + \frac{1}{2n\alpha_K} \|x_{K+1} - x_K\|^2 \\
= & \frac{\alpha_K}{2n} \left(\left\| \sum_{i=1}^n d_i(x_K) - \sum_{i=1}^n \frac{\alpha_{K,i}}{\alpha_K} d_i(\hat{z}_{K,i-1}) \right\|^2 - \left\| \sum_{i=1}^n d_i(x_K) \right\|^2 - \left\| \sum_{i=1}^n \frac{\alpha_{K,i}}{\alpha_K} d_i(\hat{z}_{K,i-1}) \right\|^2 \right) \\
& + \frac{1}{2n\alpha_K} \|x_{K+1} - x_K\|^2 \\
= & -\frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 + \frac{\alpha_K}{2n} \left\| \sum_{i=1}^n d_i(x_K) - \sum_{i=1}^n \frac{\alpha_{K,i}}{\alpha_K} d_i(\hat{z}_{K,i-1}) \right\|^2 \\
\leq & -\frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 + \frac{\alpha_K}{n} \left(\left\| \sum_{i=1}^n d_i(x_K) - \sum_{i=1}^n d_i(\hat{z}_{K,i-1}) \right\|^2 + \left\| \sum_{i=1}^n d_i(\hat{z}_{K,i-1}) - \sum_{i=1}^n \frac{\alpha_{K,i}}{\alpha_K} d_i(\hat{z}_{K,i-1}) \right\|^2 \right)
\end{aligned} \tag{21}$$

where the first two equalities are properties of the scalar product, the third equality uses (20) to drop canceling terms and the last inequality uses $\|a + b\|^2 \leq 2(\|a\|^2 + \|b\|^2)$. We bound each term separately, first,

$$\begin{aligned}
\left\| \sum_{i=1}^n d_i(x_K) - \sum_{i=1}^n d_i(\hat{z}_{K,i-1}) \right\|^2 & \leq \left(\sum_{i=1}^n \|d_i(x_K) - d_i(\hat{z}_{K,i-1})\| \right)^2 \\
& \leq \left(\sum_{i=1}^n L_i \|x_K - \hat{z}_{K,i-1}\| \right)^2 \\
& \leq \max_{i=1, \dots, n} \|x_K - \hat{z}_{K,i-1}\|^2 \left(\sum_{i=1}^n L_i \right)^2 \\
& \leq L^2 n^3 \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2.
\end{aligned} \tag{22}$$

where the first step uses the triangle inequality, the second step uses L_i Lipschicity of d_i , the third step is Hölder inequality, and the fourth step uses Claim 1. Furthermore, we

have using the triangle inequality and Cauchy-Schwartz inequality,

$$\begin{aligned}
\left\| \sum_{i=1}^n d_i(\hat{z}_{K,i-1}) - \sum_{i=1}^n \frac{\alpha_{K,i}}{\alpha_K} d_i(\hat{z}_{K,i-1}) \right\|^2 &\leq \left(\sum_{i=1}^n \left(\frac{\alpha_{K,i}}{\alpha_K} - 1 \right) \|d_i(\hat{z}_{K,i-1})\| \right)^2 \\
&\leq \sum_{i=1}^n \left(\frac{\alpha_{K,i}}{\alpha_K} - 1 \right)^2 \sum_{i=1}^n \|d_i(\hat{z}_{K,i-1})\|^2 \\
&\leq \sum_{i=1}^n \left(\frac{\alpha_{K,i}}{\alpha_K} - 1 \right)^2 \sum_{i=1}^n M_i^2 \\
&= nM^2 \sum_{i=1}^n \left(\frac{\alpha_{K,i}}{\alpha_K} - 1 \right)^2
\end{aligned} \tag{23}$$

Combining (21), (22) and (23), we obtain,

$$\begin{aligned}
&\langle \nabla F(x_K), x_{K+1} - x_K \rangle + \frac{1}{2n\alpha_K} \|x_{K+1} - x_K\|^2 \\
&\leq -\frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 + \alpha_K L^2 n^2 \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2 + \alpha_K M^2 \sum_{i=1}^n \left(\frac{\alpha_{K,i}}{\alpha_K} - 1 \right)^2,
\end{aligned}$$

which is (19) □

Claim 6 *F has L Lipschitz gradient.*

Proof : For any x, y , we have

$$\begin{aligned}
\|\nabla F(x) - \nabla F(y)\| &= \frac{1}{n} \left\| \sum_{i=1}^n d_i(x) - d_i(y) \right\| \leq \frac{1}{n} \sum_{i=1}^n \|d_i(x) - d_i(y)\| \leq \frac{1}{n} \sum_{i=1}^n L_i \|x - y\| \\
&= L \|x - y\|
\end{aligned}$$

where we used triangle inequality and L_i Lipschicity of d_i . □

Proof of Claim 2:

Using smoothness of F in Claim 6, we have from the descent Lemma [44, Lemma 1.2.3], for all $x, y \in \mathbb{R}^p$

$$F(y) \leq F(x) + \langle \nabla F(x), y - x \rangle + \frac{L}{2} \|y - x\|^2. \tag{24}$$

Choosing $y = x_{K+1}$ and $x = x_K$ in (24), using Claim 5 and Claim 1, we obtain

$$\begin{aligned}
F(x_{K+1}) &\leq F(x_K) + \langle \nabla F(x_K), x_{K+1} - x_K \rangle + \frac{L}{2} \|x_{K+1} - x_K\|^2 \\
&\leq F(x_K) - \frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 + \alpha_K L^2 n^2 \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2 + \alpha_K M^2 \sum_{i=1}^n \left(\frac{\alpha_{K,i}}{\alpha_K} - 1 \right)^2 \\
&\quad + \left(\frac{L}{2} - \frac{1}{2n\alpha_K} \right) \|x_{K+1} - x_K\|^2 \\
&\leq F(x_K) - \frac{n\alpha_K}{2} \|\nabla F(x_K)\|^2 + \left(\alpha_K L^2 n^2 + \left(\frac{Ln}{2} - \frac{1}{2\alpha_K} \right)_+ \right) \sum_{j=1}^n \alpha_{K,j}^2 \|d_j(\hat{z}_{K,j-1})\|^2 \\
&\quad + \alpha_K M^2 \sum_{i=1}^n \left(\frac{\alpha_{K,i}}{\alpha_K} - 1 \right)^2,
\end{aligned}$$

where the last inequality is obtained by Lemma 2. Since $\alpha_{K,i} \leq \alpha_K$ for all $K \in \mathbb{N}$ and $i = 1 \dots, n$, we have $0 \leq \alpha_{K,i}/\alpha_K \leq 1$, and using $(t-1)^2 \leq 1-t^2$ for all $t \in [0, 1]$

$$\left(\frac{\alpha_{K,i}}{\alpha_K} - 1 \right)^2 \leq 1 - \frac{\alpha_{K,i}^2}{\alpha_K^2} \leq 1 - \frac{\alpha_{K,i}^3}{\alpha_K^3},$$

and the result follows. \square

B Proofs for the nonsmooth setting

Proof of Theorem 1: Fix $T > 0$, we consider the sequence of functions, for each $k \in \mathbb{N}$

$$\begin{aligned}
\mathbf{w}_k &: [0, T] \mapsto \mathbb{R}^p \\
& t \mapsto \mathbf{w}(\tau_k + t)
\end{aligned}$$

From Assumption 4 and Definition 5, it is clear that all functions in the sequence are M Lipschitz. Since the sequence $(x_k)_{k \in \mathbb{N}}$ is bounded, $(\mathbf{w}_k)_{k \in \mathbb{N}}$ is also uniformly bounded, hence by Arzelà-Ascoli theorem [52, Chapter 10, Lemma 2], there is a subsequence converging uniformly, let $\mathbf{z}: [0, T] \mapsto \mathbb{R}^p$ be any such uniform limit. By discarding terms, we actually have $\mathbf{w}_k \rightarrow \mathbf{z}$ as $k \rightarrow \infty$, uniformly on $[0, T]$. Note that we have for all $t \in [0, 1]$, and all $\gamma > 0$

$$D^\gamma(\mathbf{w}_k(t)) \subset D^{\gamma + \|\mathbf{w}_k - \mathbf{z}\|_\infty}(\mathbf{z}(t)). \quad (25)$$

For all $k \in \mathbb{N}$, we set $\mathbf{v}_k \in L^2([0, T], \mathbb{R}^p)$ such that $\mathbf{v}_k = \mathbf{w}'_k$ at points where \mathbf{w}_k is differentiable (almost everywhere since it is piecewise affine). We have for all $k \in \mathbb{N}$ and all $s \in [0, T]$

$$\mathbf{w}_k(s) - \mathbf{w}_k(0) = \int_{t=0}^{t=s} \mathbf{v}_k(t) dt, \quad (26)$$

and from Definition 5, we have for almost all $t \in [0, T]$,

$$\mathbf{v}_k(t) \in -D^{\gamma(\tau_k+t)}(\mathbf{w}_k(t)). \quad (27)$$

Hence, the functions \mathbf{v}_k are uniformly bounded thanks to Assumption 4 and hence the sequence $(\mathbf{v}_k)_{k \in \mathbb{N}}$ is bounded in $L^2([0, T], \mathbb{R}^p)$ and by Banach-Alaoglu theorem [52, Section 15.1], it has a weak cluster point. Denote by \mathbf{v} a weak limit of $(\mathbf{v}_k)_{k \in \mathbb{N}}$ in $L^2([0, T], \mathbb{R}^p)$. Discarding terms, we may assume that $\mathbf{v}_k \rightarrow \mathbf{v}$ weakly in $L^2([0, T], \mathbb{R}^p)$ as $k \rightarrow \infty$ and hence, passing to the limit in (26), for all $s \in [0, T]$,

$$\mathbf{z}(s) - \mathbf{z}(0) = \int_{t=0}^{t=s} \mathbf{v}(t) dt. \quad (28)$$

By Mazur's Lemma (see for example [27]), there exists a sequence $(N_k)_{k \in \mathbb{N}}$, with $N_k \geq k$ and a sequence $\tilde{\mathbf{v}}_{k \in \mathbb{N}}$ such that for each $k \in \mathbb{N}$, $\tilde{\mathbf{v}}_k \in \text{conv}(\mathbf{v}_k, \dots, \mathbf{v}_{N_k})$ such that $\tilde{\mathbf{v}}_k$ converges strongly in $L^2([0, T], \mathbb{R}^p)$ hence pointwise almost everywhere in $[0, T]$. Using (27) and the fact that countable intersection of full measure sets has full measure, we have for almost all $t \in [0, T]$

$$\begin{aligned} \mathbf{v}(t) &= \lim_{k \rightarrow \infty} \tilde{\mathbf{v}}_k(t) \in \lim_{k \rightarrow \infty} -\text{conv} \left(\bigcup_{j=k}^{N_k} D^{\gamma(\tau_j+t)}(\mathbf{w}_j(t)) \right) \\ &\subset \lim_{k \rightarrow \infty} -\text{conv} \left(\bigcup_{j=k}^{N_k} D^{\gamma(\tau_j+t) + \|\mathbf{w}_j - \mathbf{z}\|_\infty}(\mathbf{z}(t)) \right) \\ &= -\text{conv} \left(\frac{1}{n} \sum_{i=1}^n D_i(\mathbf{z}(t)) \right) = -D(\mathbf{z}(t)). \end{aligned}$$

where we have used (25), the fact that $\lim_{\gamma \rightarrow 0} D^\gamma = \frac{1}{n} \sum_{i=1}^n D_i$ pointwise since each D_i has closed graph and the definition of D . Using (28), this shows that for almost all $t \in [0, T]$,

$$\dot{\mathbf{z}}(t) = \mathbf{v}(t) \in -D(\mathbf{z}(t)).$$

Using [7, Theorem 4.1], this shows that \mathbf{w} is an asymptotic pseudo trajectory. \square

C Lemmas and additional proofs

Lemma 2 *Let a_1, \dots, a_m be vectors in \mathbb{R}^p , then*

$$\left\| \sum_{i=1}^m a_i \right\|^2 \leq m \sum_{i=1}^m \|a_i\|^2$$

Proof : From the triangle inequality, we have

$$\left\| \sum_{i=1}^m a_i \right\|^2 \leq \left(\sum_{i=1}^m \|a_i\| \right)^2$$

Hence it suffices to prove the claim for $p = 1$. Consider the quadratic form on \mathbb{R}^m

$$Q: x \mapsto m \sum_{i=1}^m x_i^2 - \left\| \sum_{i=1}^m x_i \right\|^2.$$

We have

$$Q(x) = m(\|x\|^2 - (x^T e)^2),$$

where $e \in \mathbb{R}^m$ has unit norm and with all entries equal to $1/\sqrt{m}$. The corresponding matrix is $m(I - ee^T)$ which is positive semidefinite. This proves the result. \square

Lemma 3 *Let $(a_k)_{k \in \mathbb{N}}$ be a sequence of positive numbers, and $b, c > 0$. Then for all $m \in \mathbb{N}$*

$$\sum_{i=0}^m \frac{a_i}{b + c \sum_{j=0}^i a_j} \leq \frac{1}{c} \log \left(1 + c \frac{\sum_{i=0}^m a_i}{b} \right)$$

Proof : We have

$$\begin{aligned} \sum_{i=0}^m \frac{a_i}{b + c \sum_{j=0}^i a_j} &= \frac{1}{c} \sum_{i=0}^m \frac{a_i}{\frac{b}{c} + \sum_{j=0}^i a_j} \\ &\leq \frac{1}{c} \log \left(1 + c \frac{\sum_{i=0}^m a_i}{b} \right) \end{aligned}$$

where the last inequality follows from Lemma 6.2 in [25]. \square