



HAL
open science

Federated Expectation Maximization with heterogeneity mitigation and variance reduction

Aymeric Dieuleveut, Gersende Fort, Eric Moulines, Geneviève Robin

► **To cite this version:**

Aymeric Dieuleveut, Gersende Fort, Eric Moulines, Geneviève Robin. Federated Expectation Maximization with heterogeneity mitigation and variance reduction. NeurIPS 2021 - 35th Conference on Neural Information Processing Systems, Dec 2021, Sydney, Australia. hal-03333516v1

HAL Id: hal-03333516

<https://hal.science/hal-03333516v1>

Submitted on 3 Sep 2021 (v1), last revised 9 Nov 2021 (v3)

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

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

Federated Expectation Maximization with heterogeneity mitigation and variance reduction

Aymeric Dieuleveut

Centre de Mathématiques Appliquées
Ecole Polytechnique, France
Institut Polytechnique de Paris
aymeric.dieuleveut@polytechnique.edu

Gersende Fort

Institut de Mathématiques de Toulouse
Université de Toulouse; CNRS
UPS, Toulouse, France
gersende.fort@math.univ-toulouse.fr

Eric Moulines

Centre de Mathématiques Appliquées
Ecole Polytechnique, France
CS Dpt, HSE University, Russian Federation
eric.moulines@polytechnique.edu

Geneviève Robin

Laboratoire de Mathématiques et Modélisation d'Évry
Université d'Évry Val d'Essonne; CNRS
Évry-Courcouronnes, France
genevieve.robin@cnrs.fr

Abstract

The Expectation Maximization (EM) algorithm is the default algorithm for inference in latent variable models. As in any other field of machine learning, applications of latent variable models to very large datasets make the use of advanced parallel and distributed architectures mandatory. This paper introduces FedEM, which is the first extension of the EM algorithm to the federated learning context. FedEM is a new communication efficient method, which handles partial participation of local devices, and is robust to heterogeneous distributions of the datasets. To alleviate the communication bottleneck, FedEM compresses appropriately defined complete data sufficient statistics. We also develop and analyze an extension of FedEM to further incorporate a variance reduction scheme. In all cases, we derive finite-time complexity bounds for smooth non-convex problems. Numerical results are presented to support our theoretical findings, as well as an application to federated missing values imputation for biodiversity monitoring.

1 Introduction

The Expectation Maximization (EM) algorithm is the most popular approach for inference in latent variable models. The EM algorithm, a special instance of the Majorize/Minimize algorithm [24], was formalized by [8] and is without doubt one of the fundamental algorithms in machine learning. Applications include among many others finite mixture analysis, latent factor models inference, and missing data imputation; see [37, 28, 25, 13] and the references therein. As in any other field of machine learning, training latent variable models on very large datasets make the use of advanced parallel and distributed architectures mandatory. Federated Learning (FL) [22, 38], which exploits

the computation power of a large number of edge devices to perform distributed machine learning, is a powerful framework to achieve this goal.

The conventional EM algorithm is not suitable for FL settings. We propose several new distributed versions of the EM algorithm supporting compressed communication. More precisely, our objective is to minimize a non-convex finite-sum smooth objective function

$$\text{Argmin}_{\theta \in \Theta} F(\theta), \quad F(\theta) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\theta) + \text{R}(\theta), \quad \Theta \subseteq \mathbb{R}^d, \quad (1)$$

where n is the number of workers/devices which are connected to a central server, and the worker $\#i$ only has access to the local loss function \mathcal{L}_i ; finally R is a penalty term which may be introduced to promote sparsity, regularity, etc. In latent variable models, $\mathcal{L}_i(\theta) = -m^{-1} \sum_{j=1}^m \log p(y_{ij}; \theta)$, where $\{y_{ij}\}_{j=1}^m$ are the m observations available for worker $\#i$, and $p(y; \theta)$ is the *incomplete* likelihood. $p(y; \theta)$ is defined by marginalizing the *complete-data* likelihood $p(y, z; \theta)$ defined as the joint probability density function of the observation y and a non-observed latent variable $z \in \mathcal{Z}$, i.e. $p(y; \theta) = \int_{\mathcal{Z}} p(y, z; \theta) \mu(dz)$ where \mathcal{Z} is the *latent space* and μ is a measure on \mathcal{Z} . We focus in this paper on the case where $p(y, z; \theta)$ belongs to a curved exponential family, given by

$$p(y, z; \theta) \stackrel{\text{def}}{=} \rho(y, z) \exp \{ \langle s(y, z), \phi(\theta) \rangle - \psi(\theta) \}; \quad (2)$$

where $s(y, z) \in \mathbb{R}^q$ is the *complete-data sufficient statistics*, $\phi : \Theta \rightarrow \mathbb{R}^q$ and $\psi : \Theta \rightarrow \mathbb{R}$, $\rho : \mathcal{Y} \times \mathcal{Z} \rightarrow \mathbb{R}^+$ are vector/scalar functions.

In absence of communication constraints, the EM algorithm is a popular method to solve (1) in the curved exponential family setting. It alternates between two steps: in the Expectation (E) step, using the current value of the iterate θ_{curr} , it computes a majorizing function $\theta \mapsto \text{Q}(\theta, \theta_{\text{curr}})$ given up to an additive constant by

$$\text{Q}(\theta, \theta_{\text{curr}}) \stackrel{\text{def}}{=} -\langle \bar{s}(\theta_{\text{curr}}), \phi(\theta) \rangle + \psi(\theta) + \text{R}(\theta) \quad \text{where} \quad \bar{s}(\theta) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \bar{s}_i(\theta); \quad (3)$$

and $\bar{s}_i(\theta)$ is the i th device conditional expectation of the complete-data sufficient statistics:

$$\bar{s}_i(\theta) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij}(\theta), \quad \bar{s}_{ij}(\theta) \stackrel{\text{def}}{=} \int_{\mathcal{Z}} s(y_{ij}, z) p(z|y_{ij}; \theta) \mu(dz), \quad (4)$$

where $p(z|y_{ij}; \theta) \stackrel{\text{def}}{=} p(y_{ij}, z; \theta) / p(y_{ij}; \theta)$. As for the M step, an updated value of θ_{curr} is computed as a minimizer of $\theta \mapsto \text{Q}(\theta, \theta_{\text{curr}})$. The majorizing function is then updated with the new θ_{curr} ; this process is iterated until convergence. The EM algorithm is most useful when for any $\theta_{\text{curr}} \in \Theta$, the function $\theta \mapsto \text{Q}(\theta, \theta_{\text{curr}})$ is a convex function of the parameter θ which is solvable in θ either explicitly or with little computational effort. A major advantage of the EM algorithm stems from its invariance under homeomorphisms, contrary to classical first-order methods: the EM updates are the same for any continuous invertible re-parametrization [23].

In the FL context, the vanilla EM algorithm is affected by three major problems: (1) the communication bottleneck, (2) data heterogeneity, and (3) partial participation of the workers.

When the number of workers is large, the cost of communication becomes overwhelming. A classical technique to alleviate this problem is to use *communication compression*. Most FL algorithms are first order methods and compression is typically applied to stochastic gradients. Yet, these methods are not appropriate to solve (1) since they do not preserve the desirable homeomorphic invariance property, thus calling for an extension of the EM algorithm to the FL setting.

Since workers are often user personal devices, the issue of data heterogeneity naturally arises. Our model in Equations (1), (3) and (4) allows the local loss functions to depend on the worker $i \in [n]^*$ and the observations y_{ij} to be independent but not necessarily identically distributed. In addition, our theoretical results deal with specific behaviors for each worker $i \in [n]^*$, see e.g., A5, 7 and 9. In the FL-EM setting, heterogeneity manifests itself by the non-equality of the *local* conditional expectations of the complete-data sufficient statistics \bar{s}_i 's; modifications to the algorithms must be performed to ensure convergence at the central server.

Finally, a subset of users are potentially inactive in each learning round, being unavailable or unwilling to participate. Thus, taking into account partial participation of the workers and its impact on the convergence rate of algorithms, is a major issue.

- **FedEM.** The main contribution of our paper is a new method called FedEM, supporting communication compression, partial participation and data heterogeneity. In this algorithm, the workers compute an estimate of the *local complete-data sufficient statistics* \bar{s}_i using a minibatch of data, apply an unbiased compression operator to a noise compensated version (using a technique inspired by [17, 15]) and send the result to the central server, which performs aggregation and the M-step (i.e. the parameter update).
- **VR-FedEM.** We improve FedEM by adding a variance reduction method inspired by the SPIDER framework [9] which has recently been extended to the EM framework [10]. For both FedEM and VR-FedEM, the central server updates the iterates through a Stochastic Approximation procedure [3, 4]. When compared to FedEM, VR-FedEM additionally performs variance reduction for each worker, progressively alleviating the variance brought by the random oracles which provide approximations of the local complete-data sufficient statistics.
- **Theoretical analysis.** EM in the curved exponential family setting converges to the roots of a function h (see e.g. Section 2). We introduce a unified theoretical framework which covers the convergence of FedEM and VR-FedEM algorithms in the non-convex case and establish convergence guarantees for finding an ϵ -stationary point (see Theorem 1 and Theorem 4). In both cases, we provide the number $K_{\text{opt}}(\epsilon)$ of optimization steps and the number $K_{\text{CE}}(\epsilon)$ of computed conditional expectations \bar{s}_{ij} 's required to reach ϵ -stationarity. These results show that in the Stochastic Approximation steps of VR-FedEM, the step sizes are independent of m , the number of observations per server. Furthermore, the computational cost in terms of $\mathcal{K}_{\text{CE}}(\epsilon)$ improves on earlier results. In this respect, VR-FedEM has the same advantages as SPIDER [9] compared to SVRG [18] and SAGA [6], or as SPIDER-EM [10] compared to sEM-vr [5] and FIEM [20, 11]. Lastly, our bounds demonstrate the robustness of FedEM and VR-FedEM to data heterogeneity.
- Finally, seen as a root finding algorithm in a quantized FL setting, VR-FedEM can be compared to VR-DIANA [17]: we show that VR-FedEM does not require the step sizes to decrease with m and provides state of the art iteration complexity to reach a precision ϵ .

Notations For two vectors $a, b \in \mathbb{R}^q$, $\langle a, b \rangle$ is the Euclidean standard scalar product, and $\|\cdot\|$ denotes the associated norm. For $r \geq 1$, $\|a\|_r$ is the ℓ_r -norm of a vector a . The Hadamard product $a \odot b$ denotes the entrywise product of the two vectors a, b . By convention, vectors are column-vectors. For a matrix A , A^\top denotes its transpose and $\|A\|_F$ is its Frobenius norm; for two matrices A, B , $\langle A, B \rangle \stackrel{\text{def}}{=} \text{Trace}(B^\top A)$. For a positive integer n , set $[n]^* \stackrel{\text{def}}{=} \{1, \dots, n\}$ and $[n] \stackrel{\text{def}}{=} \{0, \dots, n\}$. The set of non-negative integers (resp. positive) is denoted by \mathbb{N} (resp. \mathbb{N}^*). The minimum (resp. maximum) of two real numbers a, b is denoted by $a \wedge b$ (resp. $a \vee b$). We will use the Bachmann-Landau notation $a(x) = O(b(x))$ to characterize an upper bound of the growth rate of $a(x)$ as being $b(x)$.

2 FedEM: Expectation Maximization algorithms for federated learning

Recall the definition of the negative penalized (normalized) log-likelihood $F(\theta)$ from (1). Along the entire paper, we make the following assumptions A1 to A3, which define the model at hand.

A1. The parameter set $\Theta \subseteq \mathbb{R}^d$ is a convex open set. The functions $R : \Theta \rightarrow \mathbb{R}$, $\phi : \Theta \rightarrow \mathbb{R}^q$, $\psi : \Theta \rightarrow \mathbb{R}$, and $\rho(y_{ij}, \cdot) : \mathcal{Z} \rightarrow \mathbb{R}_+$, $s(y_{ij}, \cdot) : \mathcal{Z} \rightarrow \mathbb{R}^q$ for $i \in [n]^*$ and $j \in [m]^*$ are measurable functions. For any $\theta \in \Theta$ and $i \in [n]^*$, the log-likelihood is bounded as $-\infty < \mathcal{L}_i(\theta) < \infty$.

A2. For all $\theta \in \Theta$ and $i \in [n]^*$, the conditional expectation $\bar{s}_i(\theta)$ is well-defined.

A3. For any $s \in \mathbb{R}^q$, the map $s \mapsto \text{Argmin}_{\theta \in \Theta} \{\psi(\theta) + R(\theta) - \langle s, \phi(\theta) \rangle\}$ exists and is unique; the singleton is denoted by $\{\mathsf{T}(s)\}$.

The EM algorithm defines a sequence $\{\theta_k, k \geq 0\}$ that can be computed recursively as $\theta_{k+1} = \mathsf{T} \circ \bar{s}(\theta_k)$, where the map T is defined in A3 and \bar{s} is defined in (3). On the other hand, the EM algorithm can be defined through a mapping in the complete-data sufficient statistics, referred to as the *expectation space*. In this setting, the EM iteration defines a \mathbb{R}^q -valued sequence $\{\widehat{S}_k, k \geq$

0} given by $\widehat{S}_{k+1} = \bar{s} \circ T(\widehat{S}_k)$. Thus, we observe that the EM algorithm admits two equivalent representations:

$$\text{(Parameter space)} \quad \theta_{k+1} = T \circ \bar{s}(\theta_k); \quad \text{(Expectation space)} \quad \widehat{S}_{k+1} = \bar{s} \circ T(\widehat{S}_k). \quad (5)$$

In this paper, we focus on the expectation space representation; see [23] for an interesting discussion on the connection of EM and mirror descent. It has been shown in [7] that if s_* is a fixed point to the EM algorithm in the expectation space, then $\theta_* \stackrel{\text{def}}{=} T(s_*)$ is a fixed point of the EM algorithm in the parameter space, i.e., $\theta_* = T \circ \bar{s}(\theta_*)$; note that the converse is also true. Define the functions h_i and h from \mathbb{R}^q to \mathbb{R}^q by

$$h(s) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n h_i(s), \quad h_i(s) \stackrel{\text{def}}{=} \bar{s}_i \circ T(s) - s. \quad (6)$$

A key property is that the fixed points of EM in the expectation space are the roots of the *mean field* $s \mapsto h(s)$ (see (3)). Therefore, convergence of EM-based algorithms is evaluated in terms of ϵ -stationarity (see [14, 10]): for all $\epsilon > 0$, there exists a (possibly random) termination time K such that

$$\mathbb{E} \left[\|h(\widehat{S}_K)\|^2 \right] \leq \epsilon. \quad (7)$$

Another key property of EM is that it is a monotonic algorithm: each iteration leads to a decrease of the negative log-likelihood i.e. $F(\theta_{k+1}) \leq F(\theta_k)$ or, equivalently in the expectation space $F \circ T(\widehat{S}_{k+1}) \leq F \circ T(\widehat{S}_k)$ (for sequences $\{\theta_k, k \geq 0\}$ and $\{\widehat{S}_k, k \geq 0\}$ given by (5)). A4 assumes that the roots of the mean field h are the roots of the gradient of $F \circ T$ (see [7] for the same assumption when studying Stochastic EM algorithms). A5 assumes global Lipschitz properties of the functions h_i 's.

A4. The function $W \stackrel{\text{def}}{=} F \circ T : \mathbb{R}^q \rightarrow \mathbb{R}$ is continuously differentiable on \mathbb{R}^q and its gradient is globally Lipschitz with constant $L_{\dot{W}}$. Furthermore, for any $s \in \mathbb{R}^q$,

$$\nabla W(s) = -B(s)h(s)$$

where $B(s)$ is a $q \times q$ positive definite matrix. In addition, there exist $0 < v_{\min} \leq v_{\max}$ such that for any $s \in \mathbb{R}^q$, the spectrum of $B(s)$ is in $[v_{\min}, v_{\max}]$.

A5. For any $i \in [n]^*$, there exists $L_i > 0$ such that for any $s, s' \in \mathbb{R}^q$,

$$\|h_i(s) - h_i(s')\| = \|(\bar{s}_i \circ T(s) - s) - (\bar{s}_i \circ T(s') - s')\| \leq L_i \|s - s'\|.$$

A Federated EM algorithm. Our first contribution, the novel algorithm FedEM is described by [algorithm 1](#). The algorithm encompasses partial participation of the workers: at iteration $\#(k+1)$, only a subset \mathcal{A}_{k+1} of active workers participate to the training, see [line 3](#). The averaged fraction of participating workers is denoted p . Each of the active workers $\#i$ computes an *unbiased* approximation $S_{k+1,i}$ ([line 6](#)) of $\bar{s}_i \circ T(\widehat{S}_k)$; conditionally to the past (see [Section 7.2](#) for a rigorous definition), these approximations are independent. The workers then transmit to the central server a compressed information about the new sufficient statistics. A naive solution would be to compress and transmit $S_{k+1,i} - \widehat{S}_k$, but data heterogeneity between servers often prevents these local differences from vanishing at the optimum, leading to large compression errors and impairing convergence of the algorithm. Following [27], a memory $V_{k,i}$ (initialized to $h_i(\widehat{S}_0)$ at $k=0$) is introduced; and the *differences* $\Delta_{k+1,i} \stackrel{\text{def}}{=} S_{k+1,i} - \widehat{S}_k - V_{k,i}$ are compressed for $i \in \mathcal{A}_{k+1}$ ([line 7](#) and [line 9](#)). These memories are updated locally: $V_{k+1,i} = V_{k,i} + \alpha \text{Quant}(\Delta_{k+1,i})$, at [line 8](#), with $\alpha > 0$ (typically set to $1/(1+\omega)$ where ω is defined in [A6](#)). On its side, the central server releases an aggregated estimate \widehat{S}_{k+1} of the complete-data sufficient statistics by averaging the quantized difference $(np)^{-1} \sum_{i \in \mathcal{A}_{k+1}} \text{Quant}(\Delta_{k+1,i})$ and by adding V_k ([line 14](#) and [line 15](#)). Then, it updates $V_{k+1} = V_k + \alpha n^{-1} \sum_{i=1}^n \text{Quant}(\Delta_{k+1,i})$, see [line 16](#). The final step consists in solving the M-step of the EM algorithm, i.e. in computing $T(\widehat{S}_{k+1})$ (see [A3](#)).

We finally state our assumption on the compression process. We consider a large class of *unbiased* compression operators Quant satisfying a variance bound:

A6. There exists $\omega \geq 0$ such that for any $s \in \mathbb{R}^q$: $\mathbb{E}[\text{Quant}(s)] = s$, and $\mathbb{E}[\|\text{Quant}(s)\|^2] \leq (1+\omega)\|s\|^2$.

Intuitively, the stronger the compression is, the larger ω will be. Remark that if no compression is used, or equivalently for all $s \in \mathbb{R}^q$, $\text{Quant}(s) = s$, then A 6 is satisfied with $\omega = 0$. An example of quantization operator satisfying A 6 is the random dithering that can be described as the random operator $\text{Quant} : \mathbb{R}^q \rightarrow \mathbb{R}^q$, $\text{Quant}(x) = (1/s_{\text{quant}})\|x\|_r \text{sign}(x) \odot [s_{\text{quant}}(\|x\|_r + \xi)]$ where $r \geq 1$ is user-defined, ξ is a uniform random variable on $[0, 1]^q$ and $s_{\text{quant}} \in \mathbb{N}^*$ is the number of levels of roundings; see [17, 2]. This operator satisfies A6 with $\omega = s_{\text{quant}}^{-1} O(q^{1/r} + q^{1/2})$; see [17, Example 1]. Another example, namely the block- p -quantization, is provided in the supplemental (see Section 6). More generally, this assumption is valid for many compression operators, for example resulting in sparsification [see. e.g. 27].

Algorithm 1: FedEM with partial participation

Data: $k_{\max} \in \mathbb{N}^*$; for $i \in [n]^*$, $V_{0,i} \in \mathbb{R}^q$; $\widehat{S}_0 \in \mathbb{R}^q$; a positive sequence $\{\gamma_{k+1}, k \in [k_{\max} - 1]\}$; $\alpha > 0$; a coefficient $p = \mathbb{E}_{\mathcal{A} \sim \mathbb{P}_{\text{PP}}}[\text{card}(\mathcal{A})]/n$.

Result: The FedEM-PP sequence: $\{\widehat{S}_k, k \in [k_{\max}]\}$

```

1 Set  $V_0 = n^{-1} \sum_{i=1}^n V_{0,i}$ ;
2 for  $k = 0, \dots, k_{\max} - 1$  do
3   Sample  $\mathcal{A}_{k+1} \sim \mathbb{P}_{\text{PP}}$ ;
4   for  $i \in \mathcal{A}_{k+1}$  do
5     (worker #i);
6     Sample  $S_{k+1,i}$ , an approximation of  $\bar{s}_i \circ T(\widehat{S}_k)$ ;
7     Set  $\Delta_{k+1,i} = S_{k+1,i} - V_{k,i} - \widehat{S}_k$ ;
8     Set  $V_{k+1,i} = V_{k,i} + \alpha \text{Quant}(\Delta_{k+1,i})$ ;
9     Send  $\text{Quant}(\Delta_{k+1,i})$  to the central server;
10  for  $i \notin \mathcal{A}_{k+1}$  do
11    (worker #i);
12    Set  $V_{k+1,i} = V_{k,i}$  (no update);
13  (the central server);
14  Set  $H_{k+1} = V_k + (np)^{-1} \sum_{i \in \mathcal{A}_{k+1}} \text{Quant}(\Delta_{k+1,i})$ ;
15  Set  $\widehat{S}_{k+1} = \widehat{S}_k + \gamma_{k+1} H_{k+1}$ ;
16  Set  $V_{k+1} = V_k + \alpha n^{-1} \sum_{i \in \mathcal{A}_{k+1}} \text{Quant}(\Delta_{k+1,i})$ ;
17  Send  $\widehat{S}_{k+1}$  and  $T(\widehat{S}_{k+1})$  to the  $n$  workers

```

The convergence analysis is under the following assumptions on the oracle $S_{k+1,i}$: for any $i \in [n]^*$, the approximations $S_{k+1,i}$ are unbiased and their conditional variances are uniformly bounded in k . For each $k \in \mathbb{N}$, denote by \mathcal{F}_k the σ -algebra generated by $\{S_{\ell,i}, \mathcal{A}_\ell; i \in [n]^*, \ell \in [k]\}$ and including the randomness inherited from the quantization operator Quant up to iteration $\#k$.

A 7. For all $k \in \mathbb{N}$, conditional to \mathcal{F}_k , $\{S_{k+1,i}\}_{i=1}^n$ are independent. Moreover, for any $i \in [n]^*$, $\mathbb{E}[S_{k+1,i} | \mathcal{F}_k] = \bar{s}_i \circ T(\widehat{S}_k)$ and there exists $\sigma_i^2 > 0$ such that for any $k \geq 0$ $\mathbb{E}[\|S_{k+1,i} - \bar{s}_i \circ T(\widehat{S}_k)\|^2 | \mathcal{F}_k] \leq \sigma_i^2$.

A7 covers both the finite-sum setting described in the introduction, and the online setting. In the finite-sum setting, \bar{s}_i is of the form $m^{-1} \sum_{j=1}^m \bar{s}_{ij}$. In that case, $S_{k+1,i}$ is the sum over a minibatch $\mathcal{B}_{k+1,i}$ of size b sampled at random in $[m]^*$, with or without replacement and independently of the history of the algorithm: we have $S_{k+1,i} = b^{-1} \sum_{j \in \mathcal{B}_{k+1,i}} \bar{s}_{ij} \circ T(\widehat{S}_k)$. In the online setting, the oracles $S_{k+1,i}$ come from an online processing of streaming informations; in that case $S_{k+1,i}$ can be computed from a minibatch of independent examples so that the conditional variance σ_i^2 , which will be inversely proportional to the size of the minibatch, can be made arbitrarily small.

Reduction of communication complexity for FL. Reducing the communication cost between workers is a crucial aspect of the Federated Learning approach [19]. In gradient based optimization, four techniques have been used to reduce the amount of communication: (i) increasing the minibatch size and reducing the number of iterations, (ii) increasing the number of *local steps* between two communication rounds, (iii) using compression, (iv) sampling clients at each step. Here, we provide a tight analysis of strategies (i), (iii) and (iv) (sampling client is part of Partial Participation).

Regarding the interest of performing multiple iterations, as analyzed for example in [21, 26] for the classical gradient settings, we note that:

- from a theoretical standpoint, tradeoffs between larger minibatch and more local iterations are unclear [36].
- *Impossibility of performing local iterations in EM.* Performing local iterations is not possible in the EM setting: one iteration of EM is the combination of two steps E and M and the M step, which corresponds here to the use of the map T , is only performed by the central server; this

remark is a fundamental specificity of the EM framework (which is not shared by the gradient framework). In applications, we usually do not want T to be available at each local node.

- However, our work allows us to perform multiple local iterations of the E step before communicating with the central server. In [algorithm 1](#), the local statistics $S_{k+1,i}$ are general enough to cover this case; see the comment above on Assumption A7.

Finally, as we do not perform local full EM iterations, we do not face the well-identified *client-drift* challenge (in the presence of heterogeneity). Yet, we stress that combining compression and heterogeneity results in other challenges: it is known in the Gradient Descent setting (see e.g. [27, 30]), that heterogeneity strongly hinders convergence in the presence of compression. To alleviate the impact of heterogeneity, we introduce the $V_{k,i}$'s memory-variables.

Convergence analysis, full participation regime. In this paragraph, we focus on the *full-participation regime*: for all $k \in [k_{\max}]^*$, $\mathcal{A}_k = [n]^*$ (consequently $p = 1$). We now present in [Theorem 1](#) our key result, from which complexity expressions are derived. The proof is postponed to [Section 7](#).

Theorem 1. *Assume A1 to A7 and set $L^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n L_i^2$, $\sigma^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n \sigma_i^2$. Let $\{\widehat{S}_k, k \in [k_{\max}]\}$ be given by [algorithm 1](#), with $\omega > 0$, $\alpha \stackrel{\text{def}}{=} (1 + \omega)^{-1}$ and $\gamma_k = \gamma \in (0, \gamma_{\max}]$ where*

$$\gamma_{\max} \stackrel{\text{def}}{=} \frac{v_{\min}}{2L_{\dot{W}}} \wedge \frac{\sqrt{n}}{2\sqrt{2}L(1+\omega)\sqrt{\omega}}. \quad (8)$$

Denote by K the uniform random variable on $[k_{\max} - 1]$. Then, taking $V_{0,i} = \mathbf{h}_i(\widehat{S}_0)$ for all $i \in [n]^*$, we have

$$v_{\min} \left(1 - \gamma \frac{L_{\dot{W}}}{v_{\min}}\right) \mathbb{E} \left[\|\mathbf{h}(\widehat{S}_K)\|^2 \right] \leq \frac{1}{\gamma k_{\max}} \left(W(\widehat{S}_0) - \min W \right) + \gamma L_{\dot{W}} \frac{1 + 5\omega}{n} \sigma^2. \quad (9)$$

When there is no compression ($\omega = 0$ so that $\text{Quant}(s) = s$), we prove that the introduction of the random variables $V_{k,i}$'s play no role whatever $\alpha > 0$ and the choice of the $V_{0,i}$'s, and we have for any $\gamma \in (0, 2v_{\min}/L_{\dot{W}})$ (see [\(31\)](#) in the supplemental)

$$\left(1 - \gamma \frac{L_{\dot{W}}}{2v_{\min}}\right) \mathbb{E} \left[\|\mathbf{h}(\widehat{S}_K)\|^2 \right] \leq \frac{1}{\gamma k_{\max}} \left(W(\widehat{S}_0) - \min W \right) + \gamma L_{\dot{W}} \frac{\sigma^2}{n}. \quad (10)$$

Optimizing the learning rate γ , we derive the following corollary (see the proof in [Section 7](#)).

Corollary 2 (of [Theorem 1](#)). *Choose $\gamma \stackrel{\text{def}}{=} \left(\frac{(W(\widehat{S}_0) - \min W)n}{k_{\max} L_{\dot{W}} (1 + 5\omega) \sigma^2} \right)^{1/2} \wedge \gamma_{\max}$. We get*

$$\mathbb{E} \left[\|\mathbf{h}(\widehat{S}_K)\|^2 \right] \leq \frac{4}{v_{\min}} \left(\sqrt{\frac{(W(\widehat{S}_0) - \min W) L_{\dot{W}} (1 + 5\omega) \sigma^2}{n k_{\max}}} \vee \frac{(W(\widehat{S}_0) - \min W)}{\gamma_{\max} k_{\max}} \right).$$

[Theorem 1](#) and [Corollary 2](#) do not require any assumption regarding the distributional heterogeneity of workers. These results remain thus valid when workers have access to data resulting from different distributions — a widespread situation in FL frameworks. Crucially, without assumptions on the heterogeneity of workers, the convergence of a “naive” implementation of compressed distributed EM (i.e. an implementation without the variables $V_{k,i}$'s) would not converge.

Let us comment the complexity to reach an ϵ -stationary point, (see [\(7\)](#)) and more precisely how the complexity evaluated in terms of the number of optimization steps depend on ω, n, σ^2 and ϵ . Since $\mathcal{K}_{\text{Opt}}(\epsilon) = k_{\max}$, from [Corollary 2](#) we have that:

$$\mathcal{K}_{\text{opt}}(\epsilon) = O \left(\frac{(1 + \omega) \sigma^2}{n \epsilon^2} \right) \vee O \left(\frac{1}{\gamma_{\max} \epsilon} \right). \quad (11)$$

Maximal learning rate and compression. The comparison of [Theorem 1](#) with the no compression case (see [\(10\)](#)) shows that compression impacts γ_{\max} by a factor proportional to $\sqrt{n}/\omega^{3/2}$ as ω increases (similar constraints were observed in the risk optimization literature, e.g. in [17, 31]).

This highlights two different regimes depending on the ratio $\sqrt{n}/\omega^{3/2}$: if the number of workers n scales at least as ω^3 , the maximal learning rate is not impacted by compression; on the other hand, for smaller numbers of workers $n \ll \omega^3$, compression can degrade the maximal learning rate. We highlight this conclusion with a small example in the case of scalar quantization for which $\omega \sim \sqrt{q}/s_{\text{quant}}$: for $q = 10^2$ and $s_{\text{quant}} = 4$ (obtaining a compression rate of a factor 16), the maximal learning rate is almost unchanged if $n \geq 16$.

Dependency on ϵ . The complexity $\mathcal{K}_{\text{opt}}(\epsilon)$ is decomposed into two terms scaling respectively as $\sigma^2\epsilon^{-2}$ and $\gamma_{\text{max}}^{-1}\epsilon^{-1}$, the first term being dominant when $\epsilon \rightarrow 0$. This first observation highlights two different regimes, i.e. a *high noise regime* corresponding to $\gamma_{\text{max}}(1+\omega)\sigma^2/(n\epsilon^{-1}) \geq 1$ where the complexity is of order $\sigma^2\epsilon^{-2}$, and a *low noise regime* where $\gamma_{\text{max}}(1+\omega)\sigma^2/(n\epsilon^{-1}) \leq 1$ and the complexity is of order $\gamma_{\text{max}}^{-1}\epsilon^{-1}$. An extreme example of the low noise regime is the case $\sigma^2 = 0$ which occurs for example in the finite-sum setting when $\bar{s}_i = m^{-1} \sum_{j=1}^m \bar{s}_{ij}$ and the oracle $S_{k+1,i}$ is equal to $\bar{s}_i \circ T(\hat{S}_k)$.

Impact of compression for ϵ -stationarity. As mentioned above, the compression simultaneously impact both the maximal learning rate (as in Equation (8)) and the complexity (see the first term in Equation (11)). Consequently, the impact of the compression depends on the balance between ω, n, σ^2 and ϵ , and we can distinguish four different “main” regimes. In the following tabular, for each of the four situations, we summarize the *increase in complexity* $\mathcal{K}_{\text{opt}}(\epsilon)$ resulting from compression.

	Complexity regime: (Dominating term in (11))	$\frac{(1+\omega)\sigma^2}{n\epsilon^2}$	$\frac{1}{\gamma_{\text{max}}\epsilon}$
γ_{max} regime: (Dominating term in (8))	Example situation	High noise σ^2 , small ϵ	Low σ^2 (e.g., large minibatch) larger ϵ
$\frac{v_{\text{min}}}{2L_{\text{W}}}$ $\frac{\sqrt{n}}{2\sqrt{2L}(1+\omega)\sqrt{\omega}}$	large ratio n/ω^3	$\times \omega$	$\times 1$
	low ratio n/ω^3	$\times \omega$	$\times \omega^{3/2}/\sqrt{n}$

Depending on the situation, the complexity can be multiplied by a factor ranging from 1 to $\omega \vee \omega^{3/2}/\sqrt{n}$. Remark that the communication cost of each iteration is typically reduced by compression of a factor at least ω . Moreover, the benefit of compression is most significant in the *low noise regime* and when the maximal learning rate is $v_{\text{min}}/(2L_{\text{W}})$ (e.g., when n large enough). We then improve the communication cost of each iteration without increasing the optimization complexity, effectively reducing the communication budget “for free”.

FedEM with partial participation and compression. In this paragraph, we extend the previous results to the *Partial Participation* (PP) regime, in which only a fraction of the workers participate to the training at each step of the learning process. This is a key feature in the FL framework, as individuals may not always be available or willing to participate [26]. To analyze the convergence in this situation, we make the following assumption.

A8. For all $k \in [k_{\text{max}} - 1]$, $\mathcal{A}_{k+1} = \{i \in [n]^* \text{ s.t. } B_{k+1,i} = 1\}$ where the random variables $B_{k+1,i}$ for $i \in [n]^*$ and $k \in [k_{\text{max}} - 1]$ are independent Bernoulli random variables with success probability $p \in (0, 1)$.

This assumption is standard in the FL literature [32, 34, 30], and can easily be extended to worker dependent probabilities of participation [16]. Algorithm 1 is naturally defined in the partial participation regime.

Usage of the control variates $(V_{k,i})_{i \in [n]^*}$ with PP. We have $V_k = n^{-1} \sum_{i=1}^n V_{k,i}$ for all $k \geq 0$ (see Proposition 12) even when the workers are not all active at iteration $\#k$. A noteworthy point is that, upon receiving $\text{Quant}(\Delta_{k+1,i})$ for all $i \in \mathcal{A}_{k+1}$, the central server computes $H_{k+1} = V_k + (np)^{-1} \sum_{i \in \mathcal{A}_{k+1}} \text{Quant}(\Delta_{k+1,i})$ and *not* $H_{k+1} = (np)^{-1} \sum_{i \in \mathcal{A}_{k+1}} (V_{k,i} + \text{Quant}(\Delta_{k+1,i}))$. Though the later solution may appear more natural, it would actually not only require to store all values $V_{k,i}$ for $i \in [n]^*$, on the central server, but also impair convergence in the heterogeneous setting. Indeed, even in the *uncompressed* regime, in which $\text{Quant}(\Delta_{k+1,i}) = \Delta_{k+1,i}$, our algorithm differs from a naive implementation of a distributed EM: FedEM computes $H_{k+1} = V_k - (np)^{-1} \sum_{i \in \mathcal{A}_{k+1}} V_{k,i} + (np)^{-1} \sum_{i \in \mathcal{A}_{k+1}} (S_{k+1,i} - \hat{S}_k)$ while a naive distributed EM would

compute $H_{k+1}^{\text{dEM}} = (np)^{-1} \sum_{i \in \mathcal{A}_{k+1}} (\mathbf{S}_{k+1,i} - \widehat{\mathbf{S}}_k)$. Such an update H_{k+1}^{dEM} is expected not to be robust to data heterogeneity as proved in [30] for the Stochastic Gradient algorithm in the FL setting. The following theorem extends Theorem 1 to the partial participation regime. Its proof is in Section 8.

Theorem 3. *Assume A1 to A8 and set $L^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n L_i^2$, $\sigma^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n \sigma_i^2$. Let $\{\widehat{\mathbf{S}}_k, k \in [k_{\max}]\}$ be given by [algorithm 1](#), run with $\alpha \stackrel{\text{def}}{=} (1 + \omega)^{-1}$ and $\gamma_k = \gamma \in (0, \gamma_{\max}]$, where*

$$\gamma_{\max} \stackrel{\text{def}}{=} \frac{v_{\min}}{2L_{\widehat{\mathbf{W}}}} \wedge \frac{p\sqrt{n}}{2\sqrt{2}L(1+\omega)\sqrt{\omega + (1-p)(1+\omega)/p}}.$$

Denote by K the uniform random variable on $[k_{\max} - 1]$. Then, taking $V_{0,i} = \mathbf{h}_i(\widehat{\mathbf{S}}_0)$ for $i \in [n]^*$, we get

$$v_{\min} \left(1 - \gamma \frac{L_{\widehat{\mathbf{W}}}}{v_{\min}}\right) \mathbb{E} \left[\|\mathbf{h}(\widehat{\mathbf{S}}_K)\|^2 \right] \leq \frac{(\mathbf{W}(\widehat{\mathbf{S}}_0) - \min \mathbf{W})}{\gamma k_{\max}} + \gamma L_{\widehat{\mathbf{W}}} \frac{1 + 5(\omega + (1-p)(1+\omega)/p)}{n} \sigma^2.$$

The above expressions can be simplified upon noting that $\omega + (1-p)(1+\omega)/p \leq (1+\omega)/p$. When $p = 1$, Theorem 1 and Theorem 3 coincide. More generally, Theorem 3 highlights that partial participation impacts both the limit variance (which increases by a factor proportional to p^{-1}) and the maximal learning rate.

3 VR-FedEM: Federated EM algorithm with variance reduction

A novel algorithm, called VR-FedEM and described by [algorithm 2](#), is derived to additionally incorporate a variance reduction scheme in FedEM. It is described in the finite-sum setting when for all $i \in [n]^*$, $\bar{\mathbf{s}}_i \stackrel{\text{def}}{=} m^{-1} \sum_{j=1}^m \bar{\mathbf{s}}_{ij}$: at each iteration $\#(t, k + 1)$, the oracle on $\bar{\mathbf{s}}_i \circ \mathbf{T}(\widehat{\mathbf{S}}_{t,k})$ will use a minibatch $\mathcal{B}_{t,k+1,i}$ of examples sampled at random (with or without replacement) in $[m]^*$.

The algorithm is decomposed into k_{out} outer loops (indexed by t), each of them having k_{in} inner loops (indexed by k). At iteration $\#(k+1)$ of the inner loops, each worker $\#i$ updates a local statistic $\mathbf{S}_{t,k+1,i}$ based on a minibatch $\mathcal{B}_{t,k+1,i}$ of its own examples $\{\bar{\mathbf{s}}_{ij}, j \in \mathcal{B}_{t,k+1,i}\}$ (see [Line 8](#)): starting from $\widehat{\mathbf{S}}_{t,0,i} \stackrel{\text{def}}{=} m^{-1} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbf{T}(\widehat{\mathbf{S}}_{t,-1})$, $\widehat{\mathbf{S}}_{t,k+1,i}$ is defined in such a way that it approximates $m^{-1} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbf{T}(\widehat{\mathbf{S}}_{t,k})$ (see [Corollary 18](#)). Then, the worker $\#i$ sends to the central server a quantization of $\Delta_{t,k+1,i}$ (see [Line 12](#)) which can be seen as an approximation of $\alpha^{-1} \{\mathbf{h}_i(\widehat{\mathbf{S}}_{t,k}) - \mathbf{h}_i(\widehat{\mathbf{S}}_{t,k-1})\}$ upon noting that the variable $V_{t,k+1,i}$ defined by [Line 10](#) approximates $\mathbf{h}_i(\widehat{\mathbf{S}}_{t,k})$ (see [Proposition 26](#)). The central server learns the mean value $V_{t,k+1} = n^{-1} \sum_{i=1}^n V_{t,k+1,i}$ (see [Line 15](#) and [Lemma 21](#)) and, by adding the quantized quantities, defines a field $\widehat{H}_{t,k+1}$ which approximates $n^{-1} \sum_{i=1}^n \mathbf{h}_i(\widehat{\mathbf{S}}_{t,k})$ (see [Proposition 24](#)). [Line 14](#) can be seen as a Stochastic Approximation update, with learning rate $\gamma_{t,k+1}$ and mean field $s \mapsto n^{-1} \sum_{i=1}^n \mathbf{h}_i(s)$ (see [\(6\)](#) for the definition of \mathbf{h}_i).

The variance reduction is encoded in the definition of $\mathbf{S}_{t,k+1,i}$, [Line 8](#). We have $\mathbf{S}_{t,k+1,i} = \mathbf{b}^{-1} \sum_{j \in \mathcal{B}_{t,k+1,i}} \bar{\mathbf{s}}_{ij} \circ \mathbf{T}(\widehat{\mathbf{S}}_{t,k}) + \Upsilon_{t,k+1,i}$. The first term is the natural approximation of $\bar{\mathbf{s}}_i \circ \mathbf{T}(\widehat{\mathbf{S}}_{t,k})$ based on a minibatch $\mathcal{B}_{t,k+1,i}$. Conditionally to the past, $\Upsilon_{t,k+1,i}$ is correlated to the first term, is biased but its bias is canceled at the beginning of each outer loop (see [Section 9.2.2](#)): $\Upsilon_{t,k+1,i}$ defines a *control variate*. Such a variance reduction technique was first proposed in the stochastic gradient setting [[29](#), [9](#), [35](#)] and then extended to the EM setting [[10](#), [12](#)]. At the end of each outer loop, the local approximations $\mathbf{S}_{t+1,0,i}$ are initialized to the full sum $m^{-1} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbf{T}(\widehat{\mathbf{S}}_{t,k_{\text{in}}})$ (see [Line 20](#)) thus canceling the bias of $\mathbf{S}_{\cdot,i}$ (see [Proposition 17](#)).

In the case where there is a single worker and no compression is used ($n = 1, \omega = 0$), VR-FedEM reduces to SPIDER-EM, which has been shown to be rate optimal for smooth, non-convex finite-sum optimization [[10](#)]. Theorem 4 studies the FL setting ($n \geq 1$ and $\omega \geq 0$): it establishes a finite time control of convergence in expectation for VR-FedEM. The assumptions [A5](#) and [A7](#) are replaced with [A9](#).

A9. For any $i \in [n]^*$ and $j \in [m]^*$, the conditional expectations $\bar{s}_{ij}(\theta)$ whatever $\theta \in \Theta$ are well defined and there exists L_{ij} such that for any $s, s' \in \mathbb{R}^q$,

$$\|(\bar{s}_{ij} \circ \mathbb{T}(s) - s) - (\bar{s}_{ij} \circ \mathbb{T}(s') - s')\| \leq L_{ij} \|s - s'\|.$$

Theorem 4. Assume A1 to 3, A4, A6 and A9. Set $L^2 \stackrel{\text{def}}{=} n^{-1} m^{-1} \sum_{i=1}^n \sum_{j=1}^m L_{ij}^2$. Let $\{\widehat{S}_{t,k}, t \in [k_{\text{out}}]^*, k \in [k_{\text{in}} - 1]\}$ be given by algorithm 2 run with $\alpha \stackrel{\text{def}}{=} 1/(1 + \omega)$, $V_{1,0,i} = \mathbf{h}_i(\widehat{S}_{1,0})$ for any $i \in [n]^*$, $\mathbf{b} \stackrel{\text{def}}{=} \lceil \frac{k_{\text{in}}}{(1+\omega)^2} \rceil$ and

$$\gamma_{t,k} = \gamma \stackrel{\text{def}}{=} \frac{1}{2\Lambda_*} = \frac{v_{\min}}{L_{\mathbb{W}}} \left(1 + 4\sqrt{2} \frac{v_{\max}}{L_{\mathbb{W}}} \frac{L}{\sqrt{n}} (1 + \omega) \left(\omega + \frac{1 + 10\omega}{8} \right)^{1/2} \right)^{-1}. \quad (12)$$

Let (τ, K) be a uniform random variable on $[k_{\text{out}}]^* \times [k_{\text{in}} - 1]$, independent of $\{\widehat{S}_{t,k}, t \in [k_{\text{out}}]^*, k \in [k_{\text{in}}]\}$. Then, it holds

$$\mathbb{E} [\|H_{\tau, K+1}\|^2] \leq \frac{2}{v_{\min} \gamma k_{\text{in}} k_{\text{out}}} \left(\mathbb{E} [\mathbb{W}(\widehat{S}_{1,0})] - \min \mathbb{W} \right), \quad (13)$$

$$\mathbb{E} [\|\mathbf{h}(\widehat{S}_{\tau, K})\|^2] \leq 2 \left(1 + \gamma^2 \frac{L^2 (1 + \omega)^2}{n} \right) \mathbb{E} [\|H_{\tau, K+1}\|^2]. \quad (14)$$

The proof is postponed to Section 9. This result is a consequence of the more general proposition Proposition 25. We make the following comments:

1. Eq. (13) provides the convergence of $\mathbb{E} [\|H_{\tau, K+1}\|^2]$, and Eq. (14) ensures that the quantity of interest $\mathbb{E} [\|\mathbf{h}(\widehat{S}_{\tau, K})\|^2]$ is controlled by $\mathbb{E} [\|H_{\tau, K+1}\|^2]$. We observe that $2 \left(1 + \gamma^2 \frac{L^2 (1 + \omega)^2}{n} \right)$ is uniformly bounded w.r.t. ω as $\gamma^2 = O_{\omega \rightarrow \infty}(\omega^{-3})$.
2. Up to our knowledge, this is the first result on Federated EM, that leverages advanced variance reduction techniques, while being robust to distribution heterogeneity (the theorem is valid without any assumption on heterogeneity) and while reducing the communication cost between workers.
3. Without compression ($\omega = 0$) and in the single-worker case ($n = 1$), Fort et al. [10] use $k_{\text{in}} = \mathbf{b}$. We observe that we recover this result as a particular case and that the recommended minibatch size \mathbf{b} decreases as $(1 + \omega)^2$.

Algorithm 2: VR-FedEM

Data: $k_{\text{out}}, k_{\text{in}}, \mathbf{b} \in \mathbb{N}^*$; for $i \in [n]^*$, $V_{1,0,i} \in \mathbb{R}^q$;
 $\widehat{S}_{\text{init}} \in \mathbb{R}^q$; a positive sequence
 $\{\gamma_{t,k+1}, t \in [k_{\text{out}}]^*, k \in [k_{\text{in}} - 1]\}$; $\alpha > 0$

Result: sequence: $\{\widehat{S}_{t,k}, t \in [k_{\text{out}}]^*, k \in [k_{\text{in}}]\}$

```

1  $\widehat{S}_{1,0} = \widehat{S}_{1,-1} = \widehat{S}_{\text{init}}, V_{1,0} = n^{-1} \sum_{i=1}^n V_{1,0,i}$ ;
2 for  $i = 1, \dots, n$  do
3    $S_{1,0,i} = \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{\text{init}})$ 
4 for  $t = 1, \dots, k_{\text{out}}$  do
5   for  $k = 0, \dots, k_{\text{in}} - 1$  do
6     for  $i = 1, \dots, n$  (worker # $i$ , locally) do
7       Sample at random a batch  $\mathcal{B}_{t,k+1,i}$  of size  $\mathbf{b}$  in  $[m]^*$ ;
8       Set  $S_{t,k+1,i} = S_{t,k,i} +$ 
9          $\mathbf{b}^{-1} \sum_{j \in \mathcal{B}_{t,k+1,i}} \left( \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) - \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k-1}) \right)$ ;
10      Set  $\Delta_{t,k+1,i} = S_{t,k+1,i} - \widehat{S}_{t,k} - V_{t,k,i}$ ;
11      Set  $V_{t,k+1,i} = V_{t,k,i} + \alpha \text{Quant}(\Delta_{t,k+1,i})$ ;
12      Send  $\text{Quant}(\Delta_{t,k+1,i})$  to the central server;
13      (the central server);
14      Set  $H_{t,k+1} = V_{t,k} + n^{-1} \sum_{i=1}^n \text{Quant}(\Delta_{t,k+1,i})$ ;
15      Set  $\widehat{S}_{t,k+1} = \widehat{S}_{t,k} + \gamma_{t,k+1} H_{t,k+1}$ ;
16      Set  $V_{t,k+1} = V_{t,k} + \alpha n^{-1} \sum_{i=1}^n \text{Quant}(\Delta_{t,k+1,i})$ ;
17      Send  $\widehat{S}_{t,k+1}$  and  $\mathbb{T}(\widehat{S}_{t,k+1})$  to the  $n$  workers;
18  $\widehat{S}_{t+1,0} = \widehat{S}_{t+1,-1} = \widehat{S}_{t,k_{\text{in}}}$ ;
19  $V_{t+1,0} = V_{t,k_{\text{in}}}$ ;
20 for  $i = 1, \dots, n$  do
21    $S_{t+1,0,i} = \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t+1,0})$ ;
22    $V_{t+1,0,i} = V_{t,k_{\text{in}},i}$ 

```

Convergence rate and optimization complexity. Our step-size γ is chosen constant and independent of $k_{\text{in}}, k_{\text{out}}$. Indeed, contrary to Theorem 1, there is no Bias-Variance trade-off (as typically observed with variance reduced methods), and the optimal choice of γ is the largest one to

ensure convergence. Consequently, since the number of optimization steps is $k_{\text{out}}k_{\text{in}}$, we have $\mathcal{K}_{\text{opt}}(\epsilon) = O(\frac{1}{\gamma\epsilon})$ where the O hides dependency on the initial value $(\mathbb{E}[\mathbb{W}(\widehat{S}_{1,0})] - \min W)$ and v_{\min} .

Impact of compression on the learning rate and ϵ -stationarity. The compression constant ω does not directly appear in Equation (13), but impacts the value of γ . Two different regimes appear:

1. if $4\sqrt{2}\frac{v_{\max}}{L_{\widehat{W}}}\frac{L}{\sqrt{n}}(1+\omega)\left(\omega + \frac{1+10\omega}{8}\right)^{1/2} \ll 1$ (i.e. we focus on the large ω, n asymptotics when $\omega^3 \ll n$), then $\gamma \simeq \frac{v_{\min}}{L_{\widehat{W}}}$ has nearly the same value as without compression [10]. The complexity is then similar to the one of SPIDER-EM [10], with a smaller communication cost. The gain from compression is maximal in this regime.
2. if $4\sqrt{2}\frac{v_{\max}}{L_{\widehat{W}}}\frac{L}{\sqrt{n}}(1+\omega)\left(\omega + \frac{1+10\omega}{8}\right)^{1/2} \gg 1$ (i.e. we focus on the large ω, n asymptotics when $\omega^3 \gg n$), then $\gamma = O\left(\frac{v_{\min}\sqrt{n}}{v_{\max}L\omega^{3/2}}\right)$ is strictly smaller than without compression. The optimization complexity is then higher to the one of SPIDER-EM¹ (by a factor proportional to $\omega^{3/2}/\sqrt{n}$) with a smaller communication cost (typically at least ω times less bits exchanged per iteration). The overall trade-off thus depends on the comparison between ω and n .

We summarize these two regimes in the following tabular, focusing on the large n , large ω asymptotic regimes. For the two regimes, we indicate the *increase in complexity* $\mathcal{K}_{\text{opt}}(\epsilon)$ resulting from compression.

	Complexity :	$\frac{1}{\gamma\epsilon}$
γ regime: (Dominating term in (12))	Example situation	
$\frac{v_{\min}}{L_{\widehat{W}}}$	large ratio n/ω^3	$\times 1$
$\frac{v_{\min}\sqrt{n}}{v_{\max}L\omega^{3/2}}$	low ratio n/ω^3	$\times \omega^{3/2}/\sqrt{n}$

Computed conditional expectations complexity \mathcal{K}_{CE} . The number of calls to conditional expectations (i.e., computing \bar{s}_{ij}) to perform k_{out} outer steps of [algorithm 2](#), each composed of k_{in} inner iterations, with n workers and mini-batches of size b is $nmk_{\text{out}} + n(2b)k_{\text{in}}k_{\text{out}} = nk_{\text{in}}k_{\text{out}}\left(\frac{m}{k_{\text{in}}} + 2b\right)$: it corresponds to one full pass on the data at the beginning of each outer iteration and two batches of size b on each worker $i \in [n]^*$, at each iteration $t, k \in [k_{\text{out}}]^* \times [k_{\text{in}} - 1]$. In order to reach an accuracy ϵ , we need $(k_{\text{in}}k_{\text{out}}\gamma)^{-1} = O(\epsilon)$ with the parameter choices in [Theorem 4](#) (esp. on b) we thus have $\mathcal{K}_{\text{CE}}(\epsilon) = O\left(\frac{n}{\epsilon\gamma}\left(\frac{m}{k_{\text{in}}} + 2\frac{k_{\text{in}}}{(1+\omega)^2}\right)\right)$. We remark that this complexity is minimized with $k_{\text{in}} = (1+\omega)\sqrt{m/2}$. We then obtain an overall complexity \mathcal{K}_{CE} of $O\left(\frac{\sqrt{m}}{\epsilon\gamma}\frac{n}{(1+\omega)}\right)$. We stress the following two points:

1. **Dependency w.r.t. m :** the complexity increases as \sqrt{m} . For $n = 1, \omega = 0$, this yields a scaling equal to $\sqrt{m}\epsilon^{-1}$ that corresponds to the optimal \mathcal{K}_{CE} of SPIDER-EM [10];
2. **Dependency w.r.t. ω .** Again, the dependency on ω depends on the regime for γ . In the (worst case regime), $\gamma = O(\sqrt{n}/\omega^{3/2})$, we get $\mathcal{K}_{\text{CE}}(\epsilon) = O_{\epsilon \rightarrow 0, \omega, n \rightarrow \infty}\left(\frac{\sqrt{m}\sqrt{n}\sqrt{\omega}}{\epsilon}\right)$, which corresponds to a sublinear increase w.r.t. ω (that compares to a linear increase in the cost of each communication).

4 Numerical illustrations

In this section, we illustrate the performance of FedEM and VR-FedEM applied to inference in Gaussian Mixture Models (GMM), on a synthetic data set and on the MNIST data set. We also present

¹As a corollary of [10, Theorem 2], the optimization complexity of SPIDER-EM is $k_{\text{out}} + k_{\text{in}}k_{\text{out}}$ that is ϵ^{-1} in order to reach ϵ -stationarity.

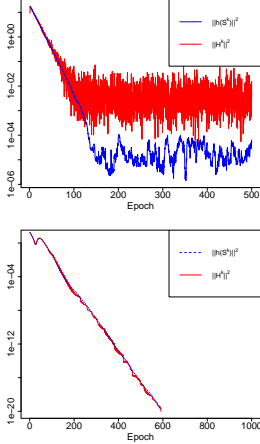


Figure 1: Trajectory of FedEM vs the number of epochs (top; blue line: $\|\mathfrak{h}(\widehat{S}^k)\|^2$; red line: $\|H_k\|^2$) and of VR-FedEM (bottom; dashed blue line: $\|\mathfrak{h}(\widehat{S}^k)\|^2$; solid red line: $\|H_{t,k}\|^2$).

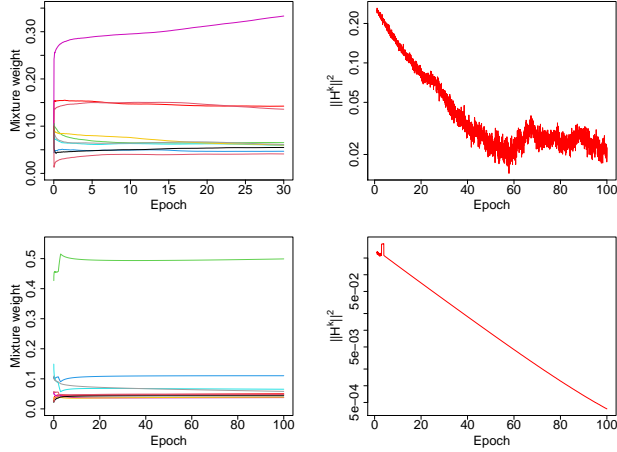


Figure 2: [Left] Evolution of the estimates of the weights π_ℓ for $\ell \in [G]^*$ by FedEM (top) and VR-FedEM (bottom) vs the number of epochs. [Right] Evolution of the squared norm of the mean field $\|H_k\|^2$ for FedEM and $\|H_{t,k}\|^2$ for VR-FedEM vs the number of epochs.

an application to Federated missing data imputation, in the context of citizen science data analysis for biodiversity monitoring with the analysis of a subsample of the eBird data set [33, 1].

Synthetic data The synthetic data are from the following GMM model: for all $\ell \in [N]^*$, and $g \in \{0, 1\}$, $\mathbb{P}(Z_\ell = g) = \pi_g$ and conditionally to $Z_\ell = g$, $Y_\ell \sim \mathcal{N}_2(\mu_g, \Sigma)$. The 2×2 covariance matrix Σ is known, and the parameters to be fitted are the weights (π_0, π_1) and the expectations (μ_0, μ_1) . The total number of examples is $N = 10^4$ and the number of agents is $n = 10^2$, and the probability of participation of servers is $p = 0.75$. FedEM and VR-FedEM are run with $\gamma = 10^{-2}$, $\omega = 1$ and $\alpha = 10^{-2}$. For FedEM, we consider the finite-sum setting when $\bar{s}_i = m^{-1} \sum_{j=1}^m \bar{s}_{ij}$ with $m = 10^2$; the oracle $S_{k+1,i}$ is obtained by a sum over a minibatch of $b = 20$ examples. For VR-FedEM, we set $b = 5$ and $k_{\text{in}} = 20$. We run the two algorithms for 500 epochs (one epoch corresponds to N conditional expectation evaluations \bar{s}_{ij}). Figure 1 shows a trajectory of FedEM and VR-FedEM in terms of $\|H_k\|^2$ for FedEM (and $\|H_{t,k}\|^2$ for VR-FedEM), along with the theoretical value of the mean field $\|\mathfrak{h}(\widehat{S})^k\|^2$. The results illustrate the effect of variance reduction, and gives insight on the variability of the trajectories resulting from the two algorithms.

MNIST Data set We perform a similar experiment on the MNIST dataset to illustrate the behaviour of FedEM and VR-FedEM on a GMM inference problem with real data. The dataset consists of $N = 7 \times 10^4$ images of handwritten digits, each with 784 pixels. We pre-process the dataset by removing 67 uninformative pixels (which are always zero across all images) to obtain $d = 717$ pixels per image. Second, we apply principal component analysis to reduce the data dimension. We keep the $d_{\text{PC}} = 20$ principal components of each observation. These N preprocessed observations are distributed at random across $n = 10^2$ servers, each containing $m = 700$ observations. We estimate a $\mathbb{R}^{d_{\text{PC}}}$ -multivariate GMM model with $G = 10$ components. Details on the multivariate Gaussian mixture model are given in the supplementary material (see section 10). Here again, \bar{s}_i is a sum over the m examples available at server $\#i$; the minibatches are independent and sampled at random in $[m]^*$ with replacement; we choose $b = 20$ and the step size is constant and set to $\gamma = 10^{-3}$. The same initial value $\widehat{S}_{\text{init}}$ is used for all experiments: we set $\widehat{S}_{\text{init}} \stackrel{\text{def}}{=} \bar{s}(\pi^0, \mu^0, \widehat{\Sigma}^0)$, where $\pi_g^0 = 1/G$ for all $g \in [G]^*$, the expectations μ_g^0 are sampled uniformly at random among the available examples, and $\widehat{\Sigma}^0$ is the empirical covariance matrix of the N examples. Figure 2 shows the sequence of parameter estimates for the weights and the squared norm of the mean field $\|H_k\|^2$ for FedEM (resp. $\|H_{t,k}\|^2$ for VR-FedEM) vs the number of epochs.

Federated missing values imputation for citizen science We develop FedMissEM, a special instance of FedEM designed to missing values imputation in the federated setting; we apply it to the analysis of part of the *eBird* data base [33, 1], a citizen science smartphone application for biodiversity monitoring. In *eBird*, citizens record wildlife observations, specifying the ecological site they visited, the date, the species and the number of observed specimens. Two major challenges occur: (i) ecological sites are visited irregularly, which leads to missing values and (ii) non-professional observers have heterogeneous wildlife counting schemes. Thus, the development of missing values prediction methods which are applicable in such federated settings is an important question.

- *Model and the FedMissEM algorithm.* I observers participate in the programme, there are J ecological sites and L time stamps. Each observer $\#i$ provides a $J \times L$ matrix X^i and a subset of indices $\Omega^i \subseteq [J]^* \times [L]^*$. For $j \in [J]^*$ and $\ell \in [L]^*$, the variable $X_{j\ell}^i$ encodes the observation that would be collected by observer $\#i$ if the site $\#j$ were visited at time stamp $\#\ell$; since there are unvisited sites, we denote by $Y^i \stackrel{\text{def}}{=} \{X_{j\ell}^i, (j, \ell) \in \Omega^i\}$ the set of observed values and $Z^i \stackrel{\text{def}}{=} \{X_{j\ell}^i, (j, \ell) \notin \Omega^i\}$ the set of unobserved values. The statistical model is parameterized by a matrix $\theta \in \mathbb{R}^{J \times L}$, where $\theta_{j\ell}$ is a scalar parameter characterizing the distribution of species individuals at site j and time stamp ℓ . For instance, $\theta_{j\ell}$ is the log-intensity of a Poisson distribution when the observations are count data or the log-odd of a binomial model when the observations are presence-absence data. This model could be extended to the case observers $\#i$ and $\#i'$ count different number of specimens on average at the same location and time stamp, for instance because they do not have access to the same material or do not have the same level of expertise: heterogeneity between observers could be modeled by using different parameters for each individual $\#i$ say $\theta^i \in \mathbb{R}^{J \times L}$. Here, we consider the case when $\theta_{j\ell}^i = \theta_{j\ell}$ for all $(j, \ell) \in [J]^* \times [L]^*$ and $i \in [I]^*$.

We further assume that the entries $\{X_{j\ell}^i, i \in [I]^*, j \in [J]^*, \ell \in [L]^*\}$ are independent with a distribution from an exponential family with respect to some reference measure ν on \mathbb{R} of the form: $x \mapsto \rho(x) \exp\{x\theta_{j\ell} - \psi(\theta_{j\ell})\}$. The function ψ is for instance defined by $\psi(\tau) = -\frac{1}{2}\tau^2$ for a Gaussian model with expectation τ and variance 1, $\psi(\tau) = \log(1 + e^\tau)$ for a Bernoulli model with success probability τ , and $\psi(\tau) = e^\tau$ for a Poisson model with intensity τ . Therefore, the joint distribution of (Y^i, Z^i) is given by

$$p_i(y^i, z^i; \theta) \stackrel{\text{def}}{=} \left(\prod_{(j, \ell) \in \Omega^i} \rho(y_{j\ell}^i) \right) \left(\prod_{(j, \ell) \notin \Omega^i} \rho(z_{j\ell}^i) \right) \exp \left(\langle s_i(y^i, z^i), \theta \rangle - \sum_{j\ell} \psi(\theta_{j\ell}) \right);$$

where $s_i(Y^i, Z^i)$ is a $J \times L$ matrix with entry $\#(j, \ell)$ given by $Y_{j\ell}^i$ if $(j, \ell) \in \Omega^i$ and $Z_{j\ell}^i$ otherwise.

In order to estimate the unknown matrix $\theta \in \mathbb{R}^{J \times L}$, we assume that θ is low-rank; we use the parameterization $\theta = UV^\top$, where $U \in \mathbb{R}^{J \times r}$ and $V \in \mathbb{R}^{L \times r}$ with $\text{rank}(\theta) = r$ and $r < \min(J, L)$. The estimator is defined as a minimizer of the negative penalized log-likelihood:

$$\min_{U \in \mathbb{R}^{J \times r}, V \in \mathbb{R}^{L \times r}} F(U, V), \quad F(U, V) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \mathcal{L}^i(UV^\top) + \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2),$$

where for $\theta \in \mathbb{R}^{J \times L}$, $\mathcal{L}^i(\theta) \stackrel{\text{def}}{=} -\log \int p_i(Y^i, z^i; \theta) \prod_{(j, \ell) \notin \Omega^i} \nu(dz_{j\ell}^i)$. Algorithm 7 in Appendix 10.2 provides the pseudo-code for FedMissEM.

- *Application to eBird data analysis.* We apply FedMissEM to the analysis of part of the *eBird* data base [33, 1] of field observations reported in France by $I = 2,465$ observers, across $J = 9,721$ sites and at $L = 525$ monthly time points. We analyze successively two data sets corresponding to observations of two relatively common species: the Common Buzzard and the Mallard. These subsamples correspond respectively to $N = 5,980$ and $N = 12,185$ field observations. The I field observers are randomly assigned into $n = 10$ groups (the observations of the field observers from the group $c \in [n]^*$ are allocated to the server $\#c$). For $c \in [n]^*$, server c contains N_c observations; in our two examples, N_c ranges between 400 and 1,500. We run the FedMissEM algorithm for 150 epochs; with $\gamma = 10^{-4}$, $\alpha = 10^{-3}$, a minibatch size $b = 10^2$, a rank $r = 2$ and $\lambda = 0$; for the distribution of the variables $X_{j\ell}^i$, we use a Gaussian distribution with unknown expectation $\theta_{j\ell}$ and variance 1. We recover aggregated temporal trends at the national French level for these two bird species by summing the estimated counts across ecological sites, for each time stamp; the trends are displayed in Figure 3.

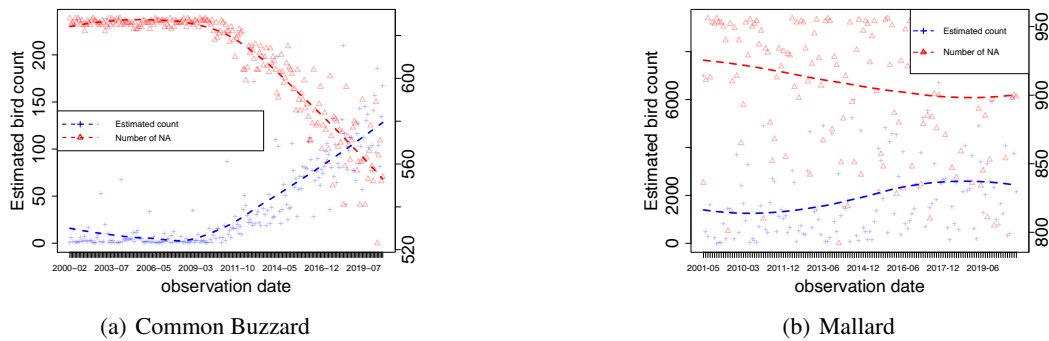


Figure 3: Estimated temporal trends for Common Buzzard (Left) and Mallard (right). Blue crosses: estimated monthly counts; Red triangles: number of missing values. Dotted lines: LOESS regressions for the estimated counts (blue) and the number of missing values (red).

5 Conclusions

We introduced the FedEM algorithm which is, to the best of our knowledge, the first method to implement the EM algorithm in a federated learning setting, and handles compression of exchanged information, data heterogeneity and partial participation. We further extended it to additionally incorporate a variance reduction scheme, and introduced the VR-FedEM algorithm. We derived complexity bounds which highlight the efficiency of the two algorithms, and illustrated our claims with numerical simulations, as well as an application to biodiversity monitoring data.

Acknowledgments The work of A. Dieuleveut and E. Moulines is partially supported by ANR-19-CHIA-0002-01 /chaire SCAI. The work of G. Fort is partially supported by the Fondation Simone et Cino del Duca under the project OpSiMorE.

Broader Impact of this work This work is mostly theoretical, and we believe it does not currently present any direct societal consequence. However, the methods described in this paper can be used to train machine learning models which could themselves have societal consequences. For instance, the deployment of machine learning models can suffer from gender and racial bias, or amplify existing inequalities.

References

- [1] ebird. 2017. ebird: An online database of bird distribution and abundance [web application]. ebird, cornell lab of ornithology, ithaca, new york. available: <http://www.ebird.org>. (accessed: 21 march 2020).
- [2] D. Alistarh, T. Hoefler, M. Johansson, N. Konstantinov, S. Khirirat, and C. Renggli. The convergence of sparsified gradient methods. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, pages 5973–5983. Curran Associates, Inc., 2018.
- [3] A. Benveniste, M. Métivier, and P. Priouret. *Adaptive Algorithms and Stochastic Approximations*. Springer Verlag, 1990.
- [4] V. S. Borkar. *Stochastic approximation*. Cambridge University Press, Cambridge; Hindustan Book Agency, New Delhi, 2008. A dynamical systems viewpoint.
- [5] J. Chen, J. Zhu, Y. Teh, and T. Zhang. Stochastic expectation maximization with variance reduction. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 7967–7977. 2018.
- [6] A. Defazio, F. Bach, and S. Lacoste-Julien. SAGA: A Fast Incremental Gradient Method With Support for Non-Strongly Convex Composite Objectives. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 1646–1654. Curran Associates, Inc., 2014.
- [7] B. Delyon, M. Lavielle, and E. Moulines. Convergence of a stochastic approximation version of the EM algorithm. *Ann. Statist.*, 27(1):94–128, 1999.
- [8] A. Dempster, N. Laird, and D. Rubin. Maximum Likelihood from Incomplete Data via the EM Algorithm. *J. Roy. Stat. Soc. B Met.*, 39(1):1–38, 1977.
- [9] C. Fang, C. Li, Z. Lin, and T. Zhang. SPIDER: Near-Optimal Non-Convex Optimization via Stochastic Path-Integrated Differential Estimator. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 689–699. Curran Associates, Inc., 2018.
- [10] G. Fort, E. Moulines, and H.-T. Wai. A Stochastic Path Integral Differential Estimator Expectation Maximization Algorithm. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 16972–16982. Curran Associates, Inc., 2020.
- [11] G. Fort, P. Gach, and E. Moulines. Fast Incremental Expectation Maximization for finite-sum optimization: non asymptotic convergence. *Statistics and Computing*, 2021. Accepted for publication.
- [12] G. Fort, E. Moulines, and H.-T. Wai. Geom-SPIDER-EM: Faster Variance Reduced Stochastic Expectation Maximization for Nonconvex Finite-Sum Optimization. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2021.
- [13] S. Frühwirth-Schnatter, G. Celeux, and C. P. Robert, editors. *Handbook of mixture analysis*. Chapman & Hall/CRC Handbooks of Modern Statistical Methods. CRC Press, Boca Raton, FL, 2019.
- [14] S. Ghadimi and G. Lan. Stochastic First- and Zeroth-Order Methods for Nonconvex Stochastic Programming. *SIAM J. Optimiz.*, 23(4):2341–2368, 2013.
- [15] E. Gorbunov, F. Hanzely, and P. Richtárik. A unified theory of SGD: Variance reduction, sampling, quantization and coordinate descent. In *International Conference on Artificial Intelligence and Statistics*, pages 680–690. PMLR, 2020.
- [16] S. Horváth and P. Richtárik. A better alternative to error feedback for communication-efficient distributed learning. In *International Conference on Learning Representations*, 2021.

- [17] S. Horváth, D. Kovalev, K. Mishchenko, S. Stich, and P. Richtárik. Stochastic distributed learning with gradient quantization and variance reduction. *arXiv preprint arXiv:1904.05115*, 2019.
- [18] R. Johnson and T. Zhang. Accelerating Stochastic Gradient Descent using Predictive Variance Reduction. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 315–323. Curran Associates, Inc., 2013.
- [19] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, R. G. L. D’Oliveira, S. E. Rouayheb, D. Evans, J. Gardner, Z. Garrett, A. Gascón, B. Ghazi, P. B. Gibbons, M. Gruteser, Z. Harchaoui, C. He, L. He, Z. Huo, B. Hutchinson, J. Hsu, M. Jaggi, T. Javidi, G. Joshi, M. Khodak, J. Konečný, A. Korolova, F. Koushanfar, S. Koyejo, T. Lepoint, Y. Liu, P. Mittal, M. Mohri, R. Nock, A. Özgür, R. Pagh, M. Raykova, H. Qi, D. Ramage, R. Raskar, D. Song, W. Song, S. U. Stich, Z. Sun, A. T. Suresh, F. Tramèr, P. Vepakomma, J. Wang, L. Xiong, Z. Xu, Q. Yang, F. X. Yu, H. Yu, and S. Zhao. Advances and Open Problems in Federated Learning. *arXiv:1912.04977 [cs, stat]*, Dec. 2019. arXiv: 1912.04977.
- [20] B. Karimi, H.-T. Wai, E. Moulines, and M. Lavielle. On the Global Convergence of (Fast) Incremental Expectation Maximization Methods. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 2837–2847. Curran Associates, Inc., 2019.
- [21] S. P. Karimireddy, S. Kale, M. Mohri, S. J. Reddi, S. U. Stich, and A. T. Suresh. SCAF-FOLD: Stochastic Controlled Averaging for On-Device Federated Learning. Oct. 2019. arXiv: 1910.06378.
- [22] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016.
- [23] F. Kunstner, R. Kumar, and M. Schmidt. Homeomorphic-invariance of em: Non-asymptotic convergence in kl divergence for exponential families via mirror descent. In *International Conference on Artificial Intelligence and Statistics*, pages 3295–3303. PMLR, 2021.
- [24] K. Lange. *MM Optimization Algorithms*. SIAM-Society for Industrial and Applied Mathematics, 2016.
- [25] G. McLachlan and T. Krishnan. *The EM algorithm and extensions*. Wiley series in probability and statistics. Wiley, 2008.
- [26] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Artificial Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- [27] K. Mishchenko, E. Gorbunov, M. Takáč, and P. Richtárik. Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- [28] K. Murphy and S. J. Russell. *Dynamic bayesian networks: representation, inference and learning*. 2002.
- [29] L. M. Nguyen, J. Liu, K. Scheinberg, and M. Takáč. Sarah: A novel method for machine learning problems using stochastic recursive gradient. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17*, page 2613–2621. JMLR.org, 2017.
- [30] C. Philippenko and A. Dieuleveut. Artemis: tight convergence guarantees for bidirectional compression in federated learning. Technical report, 2020.
- [31] C. Philippenko and A. Dieuleveut. Preserved central model for faster bidirectional compression in distributed settings. *arXiv preprint arXiv:2102.12528*, 2021.

- [32] F. Sattler, S. Wiedemann, K.-R. Müller, and W. Samek. Robust and Communication-Efficient Federated Learning From Non-i.i.d. Data. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–14, 2019.
- [33] B. L. Sullivan, C. L. Wood, M. J. Iliff, R. E. Bonney, D. Fink, and S. Kelling. eBird: A citizen-based bird observation network in the biological sciences. *Biological Conservation*, 142(10): 2282–2292, 2009.
- [34] H. Tang, C. Yu, X. Lian, T. Zhang, and J. Liu. DoubleSqueeze: Parallel Stochastic Gradient Descent with Double-pass Error-Compensated Compression. In *International Conference on Machine Learning*, pages 6155–6165. PMLR, May 2019. ISSN: 2640-3498.
- [35] Z. Wang, K. Ji, Y. Zhou, Y. Liang, and V. Tarokh. SpiderBoost and Momentum: Faster Stochastic Variance Reduction Algorithms. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 2406–2416. 2019.
- [36] B. Woodworth, K. Kshitij Patel, S. U. Stich, Z. Dai, B. Bullins, H. B. McMahan, O. Shamir, and N. Srebro. Is Local SGD Better than Minibatch SGD? *arXiv e-prints*, art. arXiv:2002.07839, Feb. 2020.
- [37] L. Xu and M. I. Jordan. On convergence properties of the EM algorithm for Gaussian mixtures. *Neural computation*, 8(1):129–151, 1996.
- [38] Q. Yang, Y. Liu, T. Chen, and Y. Tong. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):1–19, 2019.

Supplementary materials for “Federated Expectation Maximization with heterogeneity mitigation and variance reduction”

This supplementary material is organized as follows. Section 6 contains additional details on compression mechanisms satisfying A6, including an example of admissible quantization operator. Section 7 contains the pseudo-code for algorithm FedEM in the full participation regime case, and the proof of Theorem 1 – including necessary technical lemmas. Section 8 contains details concerning the extension to partial participation of the workers and the proof of Theorem 3. Section 9 is devoted to the proof of Theorem 4 concerning the convergence of VR-FedEM and necessary technical results. Finally, Section 10 contains additional details about the latent variable models used in the numerical section, as well as the pseudo code for FedMisEM.

Note that, in order to make our numerical results reproducible, code is also provided as supplementary material.

Notations For two vectors $a, b \in \mathbb{R}^q$, $\langle a, b \rangle$ is the Euclidean standard scalar product, and $\|\cdot\|$ denotes the associated norm. For $r \geq 1$, $\|a\|_r$ is the ℓ_r -norm of a vector a . The Hadamard product $a \odot b$ denotes the entrywise product of the two vectors a, b . By convention, vectors are column-vectors. For a matrix A , A^\top denotes its transpose and $\|A\|_F$ is its Frobenius norm. For a positive integer n , set $[n]^* \stackrel{\text{def}}{=} \{1, \dots, n\}$ and $[n] \stackrel{\text{def}}{=} \{0, \dots, n\}$. The set of non-negative integers (resp. positive) is denoted by \mathbb{N} (resp. \mathbb{N}^*). The minimum (resp. maximum) of two real numbers a, b is denoted by $a \wedge b$ (resp. $a \vee b$). We will use the Bachmann-Landau notation $a(x) = O(b(x))$ to characterize an upper bound of the growth rate of $a(x)$ as being $b(x)$.

We denote by $\mathcal{K}_p(\mu, \Sigma)$ the Gaussian distribution in \mathbb{R}^p , with expectation μ and covariance matrix Σ .

6 An example of quantization mechanisms: the block- p -quantization

In this section, we recall the definition of a common lossy data compression mechanism in FL (see, e.g. [27]), called block- p -quantization, and demonstrate that such quantizations satisfy the assumptions required to derive our theoretical results.

Block- p -quantization. Let $x \in \mathbb{R}^q$. Choose $\{q_\ell, \ell \in [m]^*\}$ a sequence of positive integers such that $\sum_{\ell=1}^m q_\ell = q$; and $p \in \mathbb{N}^*$. For $x \in \mathbb{R}^q$, we define the block partition

$$x = \begin{bmatrix} x_{(1)} \\ \dots \\ x_{(m)} \end{bmatrix}, \quad x_{(\ell)} \in \mathbb{R}^{q_\ell} \text{ for all } \ell \in [m]^*.$$

For all $\ell \in [m]^*$, set

$$\hat{X}_{(\ell)} \stackrel{\text{def}}{=} \|x_{(\ell)}\|_p \begin{bmatrix} \text{sign}(x_{(\ell),1}) \\ \dots \\ \text{sign}(x_{(\ell),q_\ell}) \end{bmatrix} \odot \begin{bmatrix} U_{\ell,1} \\ \dots \\ U_{\ell,q_\ell} \end{bmatrix} \quad U_{\ell,j} \stackrel{\text{indep}}{\sim} \mathcal{B}\left(\frac{|x_{(\ell),j}|}{\|x_{(\ell)}\|_p}\right), \quad (15)$$

where $x_{(\ell)} = (x_{(\ell),1}, \dots, x_{(\ell),q_\ell})^\top \in \mathbb{R}^{q_\ell}$ and $\mathcal{B}(u)$ denotes the Bernoulli random variable with success probability u . The block- p -quantization operator $\text{Quant} : \mathbb{R}^q \rightarrow \mathbb{R}^q$ is defined by

$$\text{Quant}(x) \stackrel{\text{def}}{=} \begin{bmatrix} \hat{X}_{(1)} \\ \dots \\ \hat{X}_{(m)} \end{bmatrix}. \quad (16)$$

The following Lemma ensures the block- p -quantization operator Quant satisfies the assumption A6 on the compression mechanism required by Theorem 1, Theorem 3 and Theorem 4.

Lemma 5. *Let $p \in \mathbb{N}^*$ and $\{q_\ell, \ell \in [m]^*\}$ be positive integers such that $\sum_{\ell=1}^m q_\ell = q$. For any $x \in \mathbb{R}^q$, we have*

$$\mathbb{E}[\text{Quant}(x)] = x, \quad \mathbb{E}[\|\text{Quant}(x) - x\|^2] = \sum_{\ell=1}^m (\|x_{(\ell)}\|_1 \|x_{(\ell)}\|_p - \|x_{(\ell)}\|^2),$$

where Quant is the block- p -quantization operator defined in (15) and (16). Thus, A6 holds. In particular, for $p = 2$, we may take $\omega = \max_{\ell \in [m]^*} (\sqrt{q_\ell} - 1)$.

Proof. We start by noticing that, for all $\ell \in [m]^*$, $(\text{Quant}(x))_{(\ell)} = \hat{X}_{(\ell)}$. Furthermore,

$$\begin{aligned} \mathbb{E} [\hat{X}_{(\ell)}] &= \|x_{(\ell)}\|_p \begin{bmatrix} \text{sign}(x_{(\ell),1}) \\ \dots \\ \text{sign}(x_{(\ell),q_\ell}) \end{bmatrix} \odot \begin{bmatrix} \mathbb{E}[U_{\ell,1}] \\ \dots \\ \mathbb{E}[U_{\ell,q_\ell}] \end{bmatrix} = \|x_{(\ell)}\|_p \begin{bmatrix} \text{sign}(x_{(\ell),1}) \\ \dots \\ \text{sign}(x_{(\ell),q_\ell}) \end{bmatrix} \odot \begin{bmatrix} \frac{|x_{(\ell),1}|}{\|x_{(\ell)}\|_p} \\ \dots \\ \frac{|x_{(\ell),q_\ell}|}{\|x_{(\ell)}\|_p} \end{bmatrix} \\ &= \begin{bmatrix} \text{sign}(x_{(\ell),1}) \\ \dots \\ \text{sign}(x_{(\ell),q_\ell}) \end{bmatrix} \odot \begin{bmatrix} |x_{(\ell),1}| \\ \vdots \\ |x_{(\ell),q_\ell}| \end{bmatrix} = \begin{bmatrix} x_{(\ell),1} \\ \vdots \\ x_{(\ell),q_\ell} \end{bmatrix} = x_{(\ell)}, \end{aligned}$$

which concludes the proof of the first statement. To prove the second statement, we write

$$\|\text{Quant}(x) - x\|^2 = \sum_{\ell=1}^m \|\hat{X}_{(\ell)} - x_{(\ell)}\|^2 = \sum_{\ell=1}^m \|x_{(\ell)}\|_p^2 \sum_{j=1}^{q_\ell} (U_{\ell,j} - \mathbb{E}[U_{\ell,j}])^2.$$

Since $U_{\ell,j}$ is a Bernoulli random variable with parameter $|x_{(\ell),j}|/\|x_{(\ell)}\|_p$, it holds that

$$\mathbb{E} [(U_{\ell,j} - \mathbb{E}[U_{\ell,j}])^2] = \frac{|x_{(\ell),j}| (\|x_{(\ell)}\|_p - |x_{(\ell),j}|)}{\|x_{(\ell)}\|_p^2}.$$

Hence

$$\begin{aligned} \mathbb{E} [\|\text{Quant}(x) - x\|^2] &= \sum_{\ell=1}^m \sum_{j=1}^{q_\ell} \{|x_{(\ell),j}| (\|x_{(\ell)}\|_p - |x_{(\ell),j}|)\} \\ &= \sum_{\ell=1}^m (\|x_{(\ell)}\|_1 \|x_{(\ell)}\|_p - \|x_{(\ell)}\|^2), \end{aligned}$$

which proves the second statement. In the particular case where $p = 2$, using the fact that $\|x_{(\ell)}\|_1 \leq \sqrt{q_\ell} \|x_{(\ell)}\|$, we obtain that

$$\mathbb{E} [\|\text{Quant}(x) - x\|^2] \leq \sum_{\ell=1}^m (\sqrt{q_\ell} - 1) \|x_{(\ell)}\|^2 \leq \max_{\ell \in [m]^*} (\sqrt{q_\ell} - 1) \|x\|^2,$$

which concludes the proof. \square

7 Convergence analysis of FedEM

This section contains all the elements to derive the convergence analysis of FedEM developed in Section 2 in the full participation regime. The analysis is organized as follows. First, Section 7.1 gives the pseudo code of the FedEM algorithm; Section 7.2 introduces rigorous definitions for filtrations and a technical Lemma, and Section 7.3 presents preliminary results. Then, the proof of Theorem 1 is given in Section 7.4 and the proof of Corollary 2 is in Section 7.5.

The assumptions A1 to A3 are assumed throughout this section.

7.1 Pseudo code of the FedEM algorithm

For the sake of completeness of the supplementary material, we start by recalling the pseudo code which defines the FedEM sequence in the **full participation regime**. It is given in [algorithm 3](#) below.

Algorithm 3: FedEM

Data: $k_{\max} \in \mathbb{N}^*$; for $i \in [n]^*$, $V_{0,i} \in \mathbb{R}^q$; $\widehat{S}_0 \in \mathbb{R}^q$; a positive sequence $\{\gamma_{k+1}, k \in [k_{\max} - 1]\}$; $\alpha > 0$

Result: The sequence: $\{\widehat{S}_k, k \in [k_{\max}]\}$

- 1 Set $V_0 = n^{-1} \sum_{i=1}^n V_{0,i}$;
 - 2 **for** $k = 0, \dots, k_{\max} - 1$ **do**
 - 3 **for** $i = 1, \dots, n$ **do**
 - 4 $(\text{worker } \#i)$;
 - 5 Sample $S_{k+1,i}$, an approximation of $\bar{s}_i \circ \mathbb{T}(\widehat{S}_k)$;
 - 6 Set $\Delta_{k+1,i} = S_{k+1,i} - V_{k,i} - \widehat{S}_k$;
 - 7 Set $V_{k+1,i} = V_{k,i} + \alpha \text{Quant}(\Delta_{k+1,i})$. Send $\text{Quant}(\Delta_{k+1,i})$ to the central server ;
 - 8 $(\text{the central server})$;
 - 9 Compute $H_{k+1} = V_k + n^{-1} \sum_{i=1}^n \text{Quant}(\Delta_{k+1,i})$;
 - 10 Set $\widehat{S}_{k+1} = \widehat{S}_k + \gamma_{k+1} H_{k+1}$;
 - 11 Set $V_{k+1} = V_k + \alpha n^{-1} \sum_{i=1}^n \text{Quant}(\Delta_{k+1,i})$;
 - 12 Send \widehat{S}_{k+1} and $\mathbb{T}(\widehat{S}_{k+1})$ to the n workers
-

7.2 Notations and technical lemma

In this section, we start by introducing the appropriate filtrations employed later on to define conditional expectations. Then, we present a technical lemma used in the main proof of [Theorem 1](#) (see [Section 7.4](#)).

Notations. For any random variable U , we denote by $\sigma(U)$ the sigma-algebra generated by U . For n sigma-algebras $\{\mathcal{F}_k, k \in [n]^*\}$, we denote by $\bigvee_{k=1}^n \mathcal{F}_k$ the sigma-algebra generated by $\{\mathcal{F}_k, k \in [n]^*\}$.

Definition of filtrations. Let us define the following filtrations. For any $i \in [n]^*$, we set

$$\mathcal{F}_{0,i} = \mathcal{F}_{0,i}^+ \stackrel{\text{def}}{=} \sigma\left(\widehat{S}_0; V_{0,i}\right) \quad \text{and} \quad \mathcal{F}_0 \stackrel{\text{def}}{=} \bigvee_{i=1}^n \mathcal{F}_{0,i}.$$

Then, for all $k \geq 0$,

- (i) $\mathcal{F}_{k+1/2,i} \stackrel{\text{def}}{=} \mathcal{F}_{k,i}^+ \vee \sigma(S_{k+1,i})$,
- (ii) $\mathcal{F}_{k+1,i} \stackrel{\text{def}}{=} \mathcal{F}_{k+1/2,i} \vee \sigma(\text{Quant}(\Delta_{k+1,i}))$,
- (iii) $\mathcal{F}_{k+1} \stackrel{\text{def}}{=} \bigvee_{i=1}^n \mathcal{F}_{k+1,i}$,
- (iv) $\mathcal{F}_{k+1,i}^+ \stackrel{\text{def}}{=} \mathcal{F}_{k+1,i} \vee \mathcal{F}_{k+1}$.

Note that, with these notations, for $k \geq 0$ and $i \in [n]^*$, the random variables of the FedEM sequence defined in [Algorithm 3](#) belong to the filtrations defined above as follows:

- (i) $\widehat{S}_k \in \mathcal{F}_{k,i}^+, \widehat{S}_k \in \mathcal{F}_k$,
- (ii) $S_{k+1,i}, \Delta_{k+1,i} \in \mathcal{F}_{k+1/2,i}$,
- (iii) $V_{k+1,i} \in \mathcal{F}_{k+1,i}$,
- (iv) $\widehat{S}_{k+1}, H_{k+1}, V_{k+1} \in \mathcal{F}_{k+1}$.

Note also that we have the following inclusions for filtrations: $\mathcal{F}_k \subset \mathcal{F}_{k,i}^+ \subset \mathcal{F}_{k+1/2,i} \subset \mathcal{F}_{k+1,i} \subset \mathcal{F}_{k+1}$ for all $i \in [n]^*$.

Elementary lemma. In the main proof of Theorem 1, we use the following elementary lemma.

Lemma 6. For any $x, y \in \mathbb{R}^q$ and for any $\alpha \in \mathbb{R}$, one has:

$$\|\alpha x + (1 - \alpha)y\|^2 = \alpha\|x\|^2 + (1 - \alpha)\|y\|^2 - \alpha(1 - \alpha)\|x - y\|^2.$$

Proof. The LHS is equal to

$$\alpha^2\|x\|^2 + (1 - \alpha)^2\|y\|^2 + 2\alpha(1 - \alpha)\langle x, y \rangle.$$

The RHS is equal to

$$\alpha\|x\|^2 + (1 - \alpha)\|y\|^2 - \alpha(1 - \alpha)(\|x\|^2 + \|y\|^2 - 2\langle x, y \rangle).$$

The proof is concluded upon noting that $\alpha - \alpha(1 - \alpha) = \alpha^2$ and $(1 - \alpha) - \alpha(1 - \alpha) = (1 - \alpha)^2$. \square

7.3 Preliminary results

In this section, we gather preliminary results on the control of the bias and variance of random variables of interest, which will be used in the main proof of Theorem 1. Namely, Proposition 8 controls the random field H_{k+1} , Proposition 10 controls the local increments $\Delta_{k+1,i}$ and Proposition 11 controls the memory term $V_{k,i}$.

7.3.1 Results on the memory terms V_k .

Proposition 7 shows that, even if the central server only receives the variation $\alpha^{-1}(V_{k+1,i} - V_{k,i})$ from each local worker $\#i$, it is able to compute $n^{-1}\sum_{i=1}^n V_{k+1,i}$ as soon as the quantity V_0 is correctly initialized.

Proposition 7. For any $k \in [k_{\max}]$, we have

$$V_k = \frac{1}{n} \sum_{i=1}^n V_{k,i}.$$

Proof. The proof is by induction on k . When $k = 0$, the property holds true by Line 1 in algorithm 3. Assume that the property holds for $k \leq k_{\text{in}} - 2$. Then by definition of V_{k+1} and by the induction assumption:

$$\begin{aligned} V_{k+1} &= V_k + \alpha \frac{1}{n} \sum_{i=1}^n \text{Quant}(\Delta_{k+1,i}) = \frac{1}{n} \sum_{i=1}^n (V_{k,i} + \alpha \text{Quant}(\Delta_{k+1,i})) \\ &= \frac{1}{n} \sum_{i=1}^n V_{k+1,i}. \end{aligned}$$

This concludes the induction. \square

7.3.2 Results on the random field H_{k+1} .

We compute in Proposition 8 the conditional expectation of H_{k+1} with respect to the appropriate filtration \mathcal{F}_k defined in Section 7.2, as well as an upper bound on its variance. These results are combined in an upper bound on the conditional expectation of the square norm $\|H_{k+1}\|^2$ in Corollary 9.

Proposition 8 shows that the stochastic field H_{k+1} is a (conditionally) unbiased estimator of $h(\widehat{S}_k)$. In the case of no compression (i.e. $\omega = 0$), the conditional variance of H_{k+1} is σ^2/n where σ^2 is the mean variance of the approximations $S_{k+1,i}$ over the n workers (see A7); when $\sup_i \sigma_i^2 < \infty$, the variance is inversely proportional to the number of workers n .

Proposition 8. Assume A6 and A7 and set $\sigma^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n \sigma_i^2$. For any $k \geq 0$,

$$\mathbb{E} [H_{k+1} | \mathcal{F}_k] = \mathfrak{h}(\widehat{S}_k), \quad (17)$$

$$\mathbb{E} [\|H_{k+1} - \mathbb{E} [H_{k+1} | \mathcal{F}_k]\|^2 | \mathcal{F}_k] \leq \frac{\omega}{n} \left(\frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k] \right) + \frac{\sigma^2}{n}. \quad (18)$$

Proof. Let $k \geq 0$. A6 guarantees

$$\begin{aligned} \mathbb{E} \left[\sum_{i=1}^n \text{Quant}(\Delta_{k+1,i}) \middle| \mathcal{F}_{k+1/2,i} \right] &= \sum_{i=1}^n \mathbb{E} [\text{Quant}(\Delta_{k+1,i}) | \mathcal{F}_{k+1/2,i}] \\ &= \sum_{i=1}^n \{S_{k+1,i} - V_{k,i} - \widehat{S}_k\}. \end{aligned} \quad (19)$$

Note also that, by A7, $\mathbb{E} [S_{k+1,i} | \mathcal{F}_{k,i}^+] = \bar{s}_i \circ \mathfrak{T}(\widehat{S}_k)$, and that $V_k \in \mathcal{F}_k$ and $\mathcal{F}_k \subset \mathcal{F}_{k,i}^+ \subset \mathcal{F}_{k+1/2,i}$ (see Section 7.2). Combined with (19) and using that $n^{-1} \sum_{i=1}^n V_{k,i} = V_k$ (see Proposition 7), this yields

$$\mathbb{E} [H_{k+1} | \mathcal{F}_k] = \mathbb{E} \left[n^{-1} \sum_{i=1}^n \text{Quant}(\Delta_{k+1,i}) \middle| \mathcal{F}_k \right] + V_k = \frac{1}{n} \sum_{i=1}^n \bar{s}_i \circ \mathfrak{T}(\widehat{S}_k) - \widehat{S}_k = \mathfrak{h}(\widehat{S}_k).$$

We now prove the second statement, and start by writing

$$\begin{aligned} H_{k+1} - \mathfrak{h}(\widehat{S}_k) &= \frac{1}{n} \sum_{i=1}^n \text{Quant}(\Delta_{k+1,i}) + V_k - \frac{1}{n} \sum_{i=1}^n \bar{s}_i \circ \mathfrak{T}(\widehat{S}_k) + \widehat{S}_k \\ &= \frac{1}{n} \sum_{i=1}^n \{ \text{Quant}(\Delta_{k+1,i}) - \mathbb{E} [\text{Quant}(\Delta_{k+1,i}) | \mathcal{F}_{k+1/2,i}] \} \\ &\quad + \frac{1}{n} \sum_{i=1}^n \{ S_{k+1,i} - \bar{s}_i \circ \mathfrak{T}(\widehat{S}_k) \}, \end{aligned}$$

where we applied (19) to obtain the last equality. Using the fact that $S_{k+1,i} - \bar{s}_i \circ \mathfrak{T}(\widehat{S}_k) \in \mathcal{F}_{k+1/2,i}$ and since, conditionally to \mathcal{F}_k , the workers are independent we have

$$\begin{aligned} \mathbb{E} [\|H_{k+1} - \mathfrak{h}(\widehat{S}_k)\|^2 | \mathcal{F}_k] &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\|\text{Quant}(\Delta_{k+1,i}) - \mathbb{E} [\text{Quant}(\Delta_{k+1,i}) | \mathcal{F}_{k+1/2,i}]\|^2 | \mathcal{F}_k] \\ &\quad + \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\|S_{k+1,i} - \bar{s}_i \circ \mathfrak{T}(\widehat{S}_k)\|^2 | \mathcal{F}_k]. \end{aligned}$$

The second terme in the RHS is upper bounded by $n^{-1}\sigma^2$ (see A7). For the first term, using A6 and since $\Delta_{k+1,i} \in \mathcal{F}_{k+1/2,i}$, for any $i \in [n]^*$ we have

$$\begin{aligned} &\mathbb{E} [\|\text{Quant}(\Delta_{k+1,i}) - \mathbb{E} [\text{Quant}(\Delta_{k+1,i}) | \mathcal{F}_{k+1/2,i}]\|^2 | \mathcal{F}_{k+1/2,i}] \\ &= \mathbb{E} [\|\text{Quant}(\Delta_{k+1,i})\|^2 | \mathcal{F}_{k+1/2,i}] - \|\Delta_{k+1,i}\|^2 \\ &\leq (1 + \omega) \|\Delta_{k+1,i}\|^2 - \|\Delta_{k+1,i}\|^2 = \omega \|\Delta_{k+1,i}\|^2, \end{aligned}$$

which concludes the proof upon conditioning with respect to \mathcal{F}_k . \square

Corollary 9 (of Proposition 8).

$$\mathbb{E} [\|H_{k+1}\|^2 | \mathcal{F}_k] \leq \|\mathfrak{h}(\widehat{S}_k)\|^2 + \frac{\omega}{n} \left(\frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k] \right) + \frac{\sigma^2}{n}.$$

7.3.3 Results on the local increments $\Delta_{k+1,i}$.

We compute in Proposition 10 an upper bound on the second conditional moment of $\Delta_{k+1,i}$, with respect to the appropriate filtration \mathcal{F}_k (see Section 7.2).

Proposition 10. *Assume A7. For any $i \in [n]^*$ and $k \in [k_{\max} - 1]$,*

$$\mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k] \leq \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 + \sigma_i^2.$$

Proof. Let $i \in [n]^*$ and $k \in [k_{\max} - 1]$. By A7, $\mathbb{E} [S_{k+1,i} - \widehat{S}_k | \mathcal{F}_{k,i}^+] = \mathbf{h}_i(\widehat{S}_k)$; in addition, $\widehat{S}_k \in \mathcal{F}_k$, $V_{k,i} \in \mathcal{F}_{k,i}^+$ and $\mathcal{F}_k \subset \mathcal{F}_{k,i}^+$. Hence, we get

$$\begin{aligned} \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_{k,i}^+] &= \mathbb{E} [\|S_{k+1,i} - V_{k,i} - \widehat{S}_k\|^2 | \mathcal{F}_{k,i}^+] \\ &= \|\mathbf{h}_i(\widehat{S}_k) - V_{k,i}\|^2 + \mathbb{E} [\|S_{k+1,i} - \widehat{S}_k - \mathbf{h}_i(\widehat{S}_k)\|^2 | \mathcal{F}_{k,i}^+] \\ &= \|\mathbf{h}_i(\widehat{S}_k) - V_{k,i}\|^2 + \mathbb{E} [\|S_{k+1,i} - \bar{s}_i \circ \mathbf{T}(\widehat{S}_k)\|^2 | \mathcal{F}_{k,i}^+] \\ &\stackrel{A7}{\leq} \|\mathbf{h}_i(\widehat{S}_k) - V_{k,i}\|^2 + \sigma_i^2. \end{aligned} \tag{20}$$

The proof is concluded upon noting that $\mathcal{F}_k \subset \mathcal{F}_{k,i}^+$, $\widehat{S}_k \in \mathcal{F}_k$ and $V_{k,i} \in \mathcal{F}_k$. \square

7.3.4 Results on the memory terms $V_{k,i}$.

Our final preliminary result is to compute in Proposition 11 an upper bound to control the conditional variance of the local memory terms $V_{k,i}$ with respect to the appropriate filtration \mathcal{F}_k (see Section 7.2).

Proposition 11. *Assume A5, A6 and A7; set $L^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n L_i^2$ and $\sigma^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n \sigma_i^2$. For any $k \geq 0$, set*

$$G_k \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2.$$

For any $k \in [k_{\max} - 1]$ and $\alpha \in (0, (1/(1 + \omega)))$, it holds that

$$\begin{aligned} \mathbb{E} [G_{k+1} | \mathcal{F}_k] &\leq \left(1 - \frac{\alpha}{2} + 2\gamma_{k+1}^2 \frac{L^2 \omega}{\alpha n}\right) G_k + 2\gamma_{k+1}^2 \frac{L^2}{\alpha} \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\quad + 2 \left(\alpha + \gamma_{k+1}^2 \frac{L^2 (1 + \omega)}{\alpha n}\right) \sigma^2. \end{aligned}$$

Proof. We start by computing an upper bound for the local conditional expectations $\mathbb{E} [\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_{k+1})\|^2 | \mathcal{F}_k]$, $i \in [n]^*$ and then derive the result of Proposition 11 by averaging over the n local workers.

Let $i \in [n]^*$; from Lemma 6, we have for any $s \in \mathbb{R}^q$

$$\begin{aligned} \|\mathbb{E} [V_{k+1,i} - s | \mathcal{F}_{k+1/2,i}]\|^2 &= \|(1 - \alpha)(V_{k,i} - s) + \alpha(S_{k+1,i} - \widehat{S}_k - s)\|^2 \\ &= (1 - \alpha)\|V_{k,i} - s\|^2 + \alpha\|S_{k+1,i} - \widehat{S}_k - s\|^2 - \alpha(1 - \alpha)\|\Delta_{k+1,i}\|^2. \end{aligned}$$

On the other hand,

$$\|V_{k+1,i} - \mathbb{E} [V_{k+1,i} | \mathcal{F}_{k+1/2,i}]\|^2 = \alpha^2 \|\text{Quant}(\Delta_{k+1,i}) - \mathbb{E} [\text{Quant}(\Delta_{k+1,i}) | \mathcal{F}_{k+1/2,i}]\|^2$$

and by A6 (see the proof of Proposition 8 for the same computation)

$$\mathbb{E} [\|V_{k+1,i} - \mathbb{E} [V_{k+1,i} | \mathcal{F}_{k+1/2,i}]\|^2 | \mathcal{F}_{k+1/2,i}] \leq \alpha^2 \omega \|\Delta_{k+1,i}\|^2.$$

Hence

$$\begin{aligned} \mathbb{E} \left[\|V_{k+1,i} - s\|^2 | \mathcal{F}_{k+1/2,i} \right] &\leq \mathbb{E} \left[\left\| V_{k+1,i} - s - \mathbb{E} [V_{k+1,i} - s | \mathcal{F}_{k+1/2,i}] \right\|^2 | \mathcal{F}_{k+1/2,i} \right] \\ &\quad + \mathbb{E} \left[\left\| \mathbb{E} [V_{k+1,i} - s | \mathcal{F}_{k+1/2,i}] \right\|^2 | \mathcal{F}_{k+1/2,i} \right] \\ &\leq (1 - \alpha) \|V_{k,i} - s\|^2 + \alpha \|S_{k+1,i} - \widehat{S}_k - s\|^2 + \alpha (\alpha(1 + \omega) - 1) \|\Delta_{k+1,i}\|^2. \end{aligned} \quad (21)$$

For any $\beta > 0$, using that $\|a + b\|^2 \leq (1 + \beta^2)\|a\|^2 + (1 + \beta^{-2})\|b\|^2$, we have

$$\begin{aligned} &\mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_{k+1})\|^2 | \mathcal{F}_k \right] \\ &\leq (1 + \beta^{-2}) \mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 | \mathcal{F}_k \right] + (1 + \beta^2) \mathbb{E} \left[\|\mathbf{h}_i(\widehat{S}_k) - \mathbf{h}_i(\widehat{S}_{k+1})\|^2 | \mathcal{F}_k \right] \\ &\stackrel{A5}{\leq} (1 + \beta^{-2}) \mathbb{E} \left[\mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 | \mathcal{F}_{k+1/2,i} \right] | \mathcal{F}_k \right] + (1 + \beta^2) L_i^2 \gamma_{k+1}^2 \mathbb{E} [\|H_{k+1}\|^2 | \mathcal{F}_k] \\ &\stackrel{(21)}{\leq} (1 + \beta^{-2}) \left((1 - \alpha) \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 \right. \\ &\quad \left. + \alpha \mathbb{E} [\|S_{k+1,i} - \widehat{S}_k - \mathbf{h}_i(\widehat{S}_k)\|^2 | \mathcal{F}_k] + \alpha (\alpha(1 + \omega) - 1) \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k] \right) \\ &\quad + (1 + \beta^2) L_i^2 \gamma_{k+1}^2 \mathbb{E} [\|H_{k+1}\|^2 | \mathcal{F}_k], \end{aligned}$$

where we have used (21) with $s = \mathbf{h}_i(\widehat{S}_k) \in \mathcal{F}_k \subset \mathcal{F}_{k+1/2,i}$. Choose $\beta > 0$ such that

$$\beta^{-2} \stackrel{\text{def}}{=} \begin{cases} \frac{\alpha}{2(1-\alpha)} & \text{if } \alpha \leq 2/3 \\ 1 & \text{if } \alpha \geq 2/3 \end{cases}$$

which implies that $(1 + \beta^{-2})(1 - \alpha) \leq 1 - \alpha/2$; note also that $1 \leq 1 + \beta^{-2} \leq 2$. By Corollary 9, we have (remember that $\alpha(1 + \omega) - 1 \leq 0$)

$$\begin{aligned} \mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_{k+1})\|^2 | \mathcal{F}_k \right] &\leq \left(1 - \frac{\alpha}{2}\right) \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 \\ &\quad + 2\alpha \mathbb{E} \left[\|S_{k+1,i} - \bar{s}_i \circ \mathbf{T}(\widehat{S}_k)\|^2 | \mathcal{F}_k \right] + \alpha (\alpha(1 + \omega) - 1) \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k] \\ &\quad + \frac{2}{\alpha} L_i^2 \gamma_{k+1}^2 \left(\frac{\omega}{n^2} \sum_{i=1}^n \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k] + \|\mathbf{h}(\widehat{S}_k)\|^2 + \frac{\sigma^2}{n} \right). \end{aligned}$$

Since $\alpha(1 + \omega) - 1 \leq 0$, using A7 and finally Proposition 10, we get:

$$\begin{aligned} \mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_{k+1})\|^2 | \mathcal{F}_k \right] &\leq \left(1 - \frac{\alpha}{2}\right) \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 + 2\alpha\sigma_i^2 \\ &\quad + 2\gamma_{k+1}^2 \frac{L_i^2}{\alpha} \frac{\omega}{n^2} \sum_{i=1}^n \|\mathbf{h}_i(\widehat{S}_k) - V_{k,i}\|^2 + 2\gamma_{k+1}^2 \frac{L_i^2}{\alpha} \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\quad + 2\gamma_{k+1}^2 \frac{L_i^2}{\alpha} \frac{1 + \omega}{n} \sigma^2. \end{aligned}$$

Overall, by averaging the previous inequality over all workers, we get:

$$\begin{aligned} \mathbb{E}[G_{k+1} | \mathcal{F}_k] &\leq \left(1 - \frac{\alpha}{2} + 2\gamma_{k+1}^2 \frac{L^2 \omega}{\alpha n}\right) G_k + 2\gamma_{k+1}^2 \frac{L^2}{\alpha} \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\quad + 2 \left(\alpha + \gamma_{k+1}^2 \frac{L^2 (1 + \omega)}{\alpha n} \right) \sigma^2. \end{aligned}$$

□

7.4 Proof of Theorem 1

Equipped with the necessary results, we now provide the main proof of Theorem 1. We proceed in three steps, as follows. First, for $k \geq 1$, we compute an upper bound on the average decrement

$\mathbb{E} \left[W(\widehat{S}_{k+1}) \middle| \mathcal{F}_k \right] - W(\widehat{S}_k)$ of the Lyapunov function W (defined in A4). Second, we introduce the maximal value of the learning rate. Third and finally, we deduce the result of Theorem 1 by computing the expectation w.r.t. a randomly chosen termination time K in $[k_{\max} - 1]$; in this step, we restrict the computations to the case the step sizes are constant ($\gamma_{k+1} = \gamma$ for any $k \geq 0$).

Step 1: Upper bound on the decrement. Let $k \geq 0$; from A4, we have

$$\begin{aligned} W(\widehat{S}_{k+1}) &\leq W(\widehat{S}_k) + \left\langle \nabla W(\widehat{S}_k), \widehat{S}_{k+1} - \widehat{S}_k \right\rangle + \frac{L\dot{W}}{2} \|\widehat{S}_{k+1} - \widehat{S}_k\|^2 \\ &\leq W(\widehat{S}_k) - \gamma_{k+1} \left\langle B(\widehat{S}_k) \mathbf{h}(\widehat{S}_k), H_{k+1} \right\rangle + \frac{L\dot{W}}{2} \gamma_{k+1}^2 \|H_{k+1}\|^2. \end{aligned} \quad (22)$$

Since $\widehat{S}_k \in \mathcal{F}_k$, by Proposition 8 and A4 we have

$$\mathbb{E} \left[\left\langle B(\widehat{S}_k) \mathbf{h}(\widehat{S}_k), H_{k+1} \right\rangle \middle| \mathcal{F}_k \right] = \left\langle B(\widehat{S}_k) \mathbf{h}(\widehat{S}_k), \mathbf{h}(\widehat{S}_k) \right\rangle \geq v_{\min} \|\mathbf{h}(\widehat{S}_k)\|^2. \quad (23)$$

Hence, combining (22) and (23), we have

$$\begin{aligned} \mathbb{E} \left[W(\widehat{S}_{k+1}) \middle| \mathcal{F}_k \right] &\leq W(\widehat{S}_k) - \gamma_{k+1} v_{\min} \|\mathbf{h}(\widehat{S}_k)\|^2 + \gamma_{k+1}^2 \frac{L\dot{W}}{2} \mathbb{E} [\|H_{k+1}\|^2 \middle| \mathcal{F}_k] \\ &\leq W(\widehat{S}_k) - \gamma_{k+1} v_{\min} \|\mathbf{h}(\widehat{S}_k)\|^2 + \gamma_{k+1}^2 \frac{L\dot{W}}{2} \mathbb{E} [\|H_{k+1} - \mathbb{E}[H_{k+1} \middle| \mathcal{F}_k]\|^2 \middle| \mathcal{F}_k] + \gamma_{k+1}^2 \frac{L\dot{W}}{2} \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\leq W(\widehat{S}_k) - \gamma_{k+1} v_{\min} \left(1 - \gamma_{k+1} \frac{L\dot{W}}{2v_{\min}} \right) \|\mathbf{h}(\widehat{S}_k)\|^2 + \gamma_{k+1}^2 \frac{L\dot{W}}{2} \mathbb{E} [\|H_{k+1} - \mathbb{E}[H_{k+1} \middle| \mathcal{F}_k]\|^2 \middle| \mathcal{F}_k]. \end{aligned}$$

Applying Proposition 8, we obtain that

$$\begin{aligned} \mathbb{E} \left[W(\widehat{S}_{k+1}) \middle| \mathcal{F}_k \right] &\leq W(\widehat{S}_k) - \gamma_{k+1} v_{\min} \left(1 - \gamma_{k+1} \frac{L\dot{W}}{2v_{\min}} \right) \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\quad + \gamma_{k+1}^2 \frac{L\dot{W}}{2} \frac{\omega}{n} \left(\frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\Delta_{k+1,i}\|^2 \middle| \mathcal{F}_k] \right) + \gamma_{k+1}^2 \frac{L\dot{W}}{2n} \sigma^2. \end{aligned} \quad (24)$$

Finally, using Proposition 10 and (24), we get:

$$\begin{aligned} \mathbb{E}[W(\widehat{S}_{k+1}) \middle| \mathcal{F}_k] &\leq W(\widehat{S}_k) - \gamma_{k+1} v_{\min} \left(1 - \gamma_{k+1} \frac{L\dot{W}}{2v_{\min}} \right) \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\quad + \gamma_{k+1}^2 \frac{L\dot{W}}{2} \frac{\omega}{n} G_k + \gamma_{k+1}^2 \frac{L\dot{W}}{2n} (1 + \omega) \sigma^2, \end{aligned} \quad (25)$$

where

$$G_k \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2.$$

Step 2: Maximal learning rate γ_{k+1} **when** $\omega \neq 0$. From Proposition 11, for any non-increasing positive sequence $\{\gamma_k, k \in [k_{\max} - 1]\}$ such that

$$\gamma_{k+1}^2 \leq \frac{\alpha^2}{8L^2} \frac{n}{\omega},$$

and for any positive sequence $\{C_k, k \in [k_{\max} - 1]\}$, it holds

$$\begin{aligned} C_{k+1} \mathbb{E} [G_{k+1} \middle| \mathcal{F}_k] &\leq C_{k+1} \left(1 - \frac{\alpha}{4} \right) G_k \\ &\quad + C_{k+1} \gamma_{k+1}^2 \frac{2}{\alpha} L^2 \|\mathbf{h}(\widehat{S}_k)\|^2 + 2C_{k+1} \left(\alpha + \gamma_{k+1}^2 \frac{L^2}{\alpha} \frac{1 + \omega}{n} \right) \sigma^2. \end{aligned} \quad (26)$$

Combining equations (25) and (26), we thus have

$$\begin{aligned} \mathbb{E}[\mathbb{W}(\widehat{S}_{k+1})|\mathcal{F}_k] + C_{k+1}\mathbb{E}[G_{k+1}|\mathcal{F}_k] &\leq \mathbb{W}(\widehat{S}_k) + C_k G_k \\ &\quad - \gamma_{k+1} v_{\min} \left(1 - \gamma_{k+1} \frac{L\dot{\mathbb{W}}}{2v_{\min}} - \frac{C_{k+1}}{v_{\min}} \gamma_{k+1} \frac{2}{\alpha} L^2 \right) \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\quad + \left(\gamma_{k+1}^2 \frac{L\dot{\mathbb{W}}}{2} \frac{\omega}{n} - C_k + C_{k+1} - C_{k+1} \frac{\alpha}{4} \right) G_k \\ &\quad + \left\{ 2\alpha C_{k+1} + \gamma_{k+1}^2 \frac{(1+\omega)}{n} \left(\frac{L\dot{\mathbb{W}}}{2} + 2C_{k+1} \frac{L^2}{\alpha} \right) \right\} \sigma^2. \end{aligned}$$

We choose the sequence $\{C_k\}$ as follows:

$$C_k \stackrel{\text{def}}{=} \gamma_k^2 \frac{2L\dot{\mathbb{W}}}{\alpha} \frac{\omega}{n};$$

the sequence satisfies $C_{k+1} \leq C_k$ (since $\gamma_{k+1} \leq \gamma_k$) and $\gamma_{k+1}^2 L\dot{\mathbb{W}}\omega/(2n) \leq C_{k+1}\alpha/4$. By convention, $\gamma_0 \in [\gamma_1, +\infty)$. Therefore

$$\mathbb{E}[\mathbb{W}(\widehat{S}_{k+1})|\mathcal{F}_k] + \gamma_{k+1}^2 \frac{2L\dot{\mathbb{W}}}{\alpha} \frac{\omega}{n} \mathbb{E}[G_{k+1}|\mathcal{F}_k] \leq \mathbb{W}(\widehat{S}_k) + \gamma_k^2 \frac{2L\dot{\mathbb{W}}}{\alpha} \frac{\omega}{n} G_k \quad (27)$$

$$- \gamma_{k+1} v_{\min} \left(1 - \gamma_{k+1} \frac{L\dot{\mathbb{W}}}{2v_{\min}} \left\{ 1 + 8\gamma_{k+1}^2 \frac{\omega}{\alpha^2 n} L^2 \right\} \right) \|\mathbf{h}(\widehat{S}_k)\|^2 \quad (28)$$

$$+ 4\gamma_{k+1}^2 L\dot{\mathbb{W}} \frac{\omega}{n} \left\{ 1 + \frac{(1+\omega)}{8\omega} \left(1 + \gamma_{k+1}^2 8 \frac{L^2 \omega}{\alpha^2 n} \right) \right\} \sigma^2. \quad (29)$$

Step 3: Computing the expectation. Let us apply the expectations, sum from $k = 0$ to $k = k_{\max} - 1$, and divide by k_{\max} . This yields

$$\begin{aligned} &\frac{v_{\min}}{k_{\max}} \sum_{k=0}^{k_{\max}-1} \gamma_{k+1} \left(1 - \gamma_{k+1} \frac{L\dot{\mathbb{W}}}{2v_{\min}} \left\{ 1 + 8\gamma_{k+1}^2 \frac{\omega}{\alpha^2 n} L^2 \right\} \right) \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\leq k_{\max}^{-1} \left\{ \mathbb{W}(\widehat{S}_0) + \gamma_0^2 \frac{2L\dot{\mathbb{W}}}{\alpha} \frac{\omega}{n} G_0 - \mathbb{E}[\mathbb{W}(\widehat{S}_{k_{\max}})] - \gamma_{k_{\max}}^2 \frac{2L\dot{\mathbb{W}}}{\alpha} \frac{\omega}{n} \mathbb{E}[G_{k_{\max}}] \right\} \\ &\quad + 4L\dot{\mathbb{W}} \frac{\omega}{n} \frac{1}{k_{\max}} \sum_{k=0}^{k_{\max}-1} \gamma_{k+1}^2 \left\{ 1 + \frac{(1+\omega)}{8\omega} \left(1 + \gamma_{k+1}^2 8 \frac{L^2 \omega}{\alpha^2 n} \right) \right\} \sigma^2. \end{aligned}$$

We now focus on the case when $\gamma_{k+1} = \gamma$ for any $k \geq 0$. Denote by K a uniform random variable on $[k_{\max} - 1]$, independent of the path $\{\widehat{S}_k, k \in [k_{\max}]\}$. Since $\gamma^2 \leq \alpha^2 n / (8L^2 \omega)$, we have

$$1 + 8\gamma^2 \frac{\omega}{\alpha^2 n} L^2 \leq 2.$$

This yields

$$\begin{aligned} &v_{\min} \gamma \left(1 - \gamma \frac{L\dot{\mathbb{W}}}{v_{\min}} \right) \mathbb{E} \left[\|\mathbf{h}(\widehat{S}_K)\|^2 \right] \\ &\leq k_{\max}^{-1} \left\{ \mathbb{W}(\widehat{S}_0) + \gamma^2 \frac{2L\dot{\mathbb{W}}}{\alpha} \frac{\omega}{n} G_0 - \mathbb{E}[\mathbb{W}(\widehat{S}_{k_{\max}})] - \gamma^2 \frac{2L\dot{\mathbb{W}}}{\alpha} \frac{\omega}{n} \mathbb{E}[G_{k_{\max}}] \right\} \\ &\quad + 4L\dot{\mathbb{W}} \frac{\omega}{n} \gamma^2 \left\{ 1 + \frac{(1+\omega)}{4\omega} \right\} \sigma^2. \end{aligned} \quad (30)$$

Note that $4(1 + (1+\omega)/(4\omega)) = (5\omega + 1)/\omega$.

Step 4. Conclusion (when $\omega \neq 0$). By choosing $V_{0,i} = \mathbf{h}_i$ for any $i \in [n]^*$, we have $G_0 = 0$. The roots of $\gamma \mapsto \gamma(1 - \gamma L\dot{\mathbb{W}}/v_{\min})$ are 0 and $v_{\min}/L\dot{\mathbb{W}}$ and its maximum is reached at $v_{\min}/(2L\dot{\mathbb{W}})$: this function is increasing on $(0, v_{\min}/(2L\dot{\mathbb{W}})]$. We therefore choose $\gamma \in (0, \gamma_{\max}(\alpha)]$ where

$$\gamma_{\max}(\alpha) \stackrel{\text{def}}{=} \min \left(\frac{v_{\min}}{2L\dot{\mathbb{W}}}; \frac{\alpha}{2\sqrt{2}L} \frac{\sqrt{n}}{\sqrt{\omega}} \right)$$

Finally, since $\alpha \in (0, 1/(1+\omega)]$, we choose $\alpha = 1/(1+\omega)$. This yields

$$\gamma_{\max} \stackrel{\text{def}}{=} \min \left(\frac{v_{\min}}{2L\dot{\mathbb{W}}}; \frac{1}{2\sqrt{2}L} \frac{\sqrt{n}}{\sqrt{\omega}(1+\omega)} \right).$$

Case $\omega = 0$. From (25), applying the expectation we have

$$\gamma_{k+1} v_{\min} \left(1 - \gamma_{k+1} \frac{L_{\dot{W}}}{2v_{\min}} \right) \mathbb{E} \left[\|\mathbf{h}(\widehat{S}_k)\|^2 \right] \leq \mathbb{E} \left[\mathbf{W}(\widehat{S}_k) \right] - \mathbb{E} \left[\mathbf{W}(\widehat{S}_{k+1}) \right] + \gamma_{k+1}^2 \frac{L_{\dot{W}} \sigma^2}{2n}.$$

We now sum from $k = 0$ to $k = k_{\max} - 1$ and then divide by k_{\max} . In the case $\gamma_{k+1} = \gamma$, we have

$$\gamma v_{\min} \left(1 - \gamma \frac{L_{\dot{W}}}{2v_{\min}} \right) \mathbb{E} \left[\|\mathbf{h}(\widehat{S}_K)\|^2 \right] \leq k_{\max}^{-1} \left(\mathbb{E} \left[\mathbf{W}(\widehat{S}_0) \right] - \min \mathbf{W} \right) + \gamma^2 \frac{L_{\dot{W}} \sigma^2}{2n}. \quad (31)$$

Remark on the maximal learning rate. The condition $\gamma_{k+1} \leq \frac{\alpha}{2\sqrt{2}L} \frac{\sqrt{n}}{\sqrt{\omega}}$ is used twice in the proof:

1. To ensure that $\left(1 - \gamma_{k+1} \frac{L_{\dot{W}}}{2v_{\min}} \left\{ 1 + 8\gamma_{k+1}^2 \frac{\omega}{\alpha^2 n} L^2 \right\} \right) \geq \left(1 - \gamma_{k+1} \frac{L_{\dot{W}}}{v_{\min}} \right)$ in order to obtain Equation (30).
2. To ensure that the process $(G_k)_{k \geq 0}$ is ‘‘pseudo-contractive’’ (i.e., satisfies a recursion of the form $u_{k+1} \leq \rho u_k + v_k$, with $\rho < 1$) in Proposition 11.

A more detailed analysis can get rid of this condition (and thus the dependency $\gamma = O_{\omega \rightarrow \infty}(\omega^{-3/2})$, as we recall that $\alpha^1 \propto_{\omega \rightarrow \infty} \omega$) for the *first point*. Indeed, we ultimately only require

$$\left(1 - \gamma_{k+1} \frac{L_{\dot{W}}}{2v_{\min}} \left\{ 1 + 8\gamma_{k+1}^2 \frac{\omega}{\alpha^2 n} L^2 \right\} \right) \geq \frac{1}{2} \quad (32)$$

to conclude the proof. This is for example satisfied if $\gamma_{k+1} \frac{L_{\dot{W}}}{2v_{\min}} \leq \frac{1}{4}$ and $8\gamma_{k+1}^3 \frac{L_{\dot{W}}}{2v_{\min}} \frac{\omega}{\alpha^2 n} L^2 \leq \frac{1}{4}$. This approach results in a better asymptotic dependency of the maximal learning rate w.r.t. ω to obtain Equation (32): $\gamma = O_{\omega \rightarrow \infty}(\omega^{-1})$. However, the condition $\gamma_{k+1} \leq \frac{\alpha}{2\sqrt{2}L} \frac{\sqrt{n}}{\sqrt{\omega}}$ seems to be *necessary* to obtain the *second point* and Proposition 11. The possibility of providing a similar result to Proposition 11 without the $\omega^{-3/2}$ dependency, is an interesting open problem.

7.5 Proof of Corollary 2

In (9), the RHS is of the form $A/\gamma + \gamma B$ for some positive constants A, B : we have $A/\gamma + \gamma B \geq 2\sqrt{AB}$ with equality reached with $\gamma_{\star} \stackrel{\text{def}}{=} \sqrt{A/B}$. Hence, we set

$$\gamma_{\star} \stackrel{\text{def}}{=} \frac{1}{\sigma} \left(\frac{n \left(\mathbf{W}(\widehat{S}_0) - \min \mathbf{W} \right)}{L_{\dot{W}}(1+5\omega)} \right)^{1/2} \frac{1}{\sqrt{k_{\max}}}.$$

If $\gamma_{\star} \leq \gamma_{\max}$, then let us apply (9) with $\gamma = \gamma_{\star}$ which yields a RHS given by $2\sqrt{A/B}$ i.e.

$$2\sigma \left(\left(\mathbf{W}(\widehat{S}_0) - \min \mathbf{W} \right) L_{\dot{W}} \frac{(1+5\omega)}{n} \right)^{1/2} \frac{1}{\sqrt{k_{\max}}}.$$

If $\gamma_{\star} \geq \gamma_{\max}$, we write

$$\frac{A}{\gamma_{\max}} + B\gamma_{\max} \leq \frac{A}{\gamma_{\max}} + \frac{A}{\gamma_{\max}} \frac{\gamma_{\max}^2 B}{A} = \frac{A}{\gamma_{\max}} + \frac{A}{\gamma_{\max}} \frac{\gamma_{\max}^2}{\gamma_{\star}^2} \leq 2 \frac{A}{\gamma_{\max}}.$$

and the RHS is upper bounded by

$$2 \frac{\mathbf{W}(\widehat{S}_0) - \min \mathbf{W}}{\gamma_{\max} k_{\max}}.$$

Finally, in the LHS of (9), we have

$$1 - \gamma \frac{L_{\dot{W}}}{v_{\min}} \geq 1 - \gamma_{\max} \frac{L_{\dot{W}}}{v_{\min}} \geq 1 - \frac{v_{\min}}{2L_{\dot{W}}} \frac{L_{\dot{W}}}{v_{\min}} = \frac{1}{2}.$$

This concludes the proof.

8 Partial Participation case

In this section, we generalize the result of Theorem 1 to the *partial participation case*. This extra scheme could be incorporated into the main proof, but we choose to present it separately to improve the readability of the main proof in Section 7. We first provide an equivalent description of [algorithm 1](#) in Section 8.1; [algorithm 4](#) will be used throughout this section. Then, we introduce a new family of filtrations. In Section 8.3, we first establish preliminary results and then give the proof of Theorem 3 in Section 8.4.

The assumptions A1 to A3 hold throughout this section.

8.1 An equivalent algorithm

In this Section, we describe an equivalent algorithm, that outputs the same result as Algorithm 1, and for which the analysis is conducted.

Algorithm 4: FedEM with partial participation

Data: $k_{\max} \in \mathbb{N}^*$; for $i \in [n]^*$, $V_{0,i} \in \mathbb{R}^a$; $\widehat{S}_0 \in \mathbb{R}^a$; a positive sequence $\{\gamma_{k+1}, k \in [k_{\max} - 1]\}$; $\alpha > 0$; $p \in (0, 1)$.

Result: The FedEM-PP sequence: $\{\widehat{S}_k, k \in [k_{\max}]\}$

- 1 Set $V_0 = n^{-1} \sum_{i=1}^n V_{0,i}$;
 - 2 **for** $k = 0, \dots, k_{\max} - 1$ **do**
 - 3 **for** $i = 1, \dots, n$ **do**
 - 4 (worker # i);
 - 5 Sample $S_{k+1,i}$, an approximation of $\bar{s}_i \circ \mathsf{T}(\widehat{S}_k)$;
 - 6 Set $\Delta_{k+1,i} = S_{k+1,i} - V_{k,i} - \widehat{S}_k$;
 - 7 Sample a Bernoulli r.v. $B_{k+1,i}$ with success probability p ;
 - 8 Set $V_{k+1,i} = V_{k,i} + \alpha B_{k+1,i} \text{Quant}(\Delta_{k+1,i})$. ;
 - 9 Send $B_{k+1,i} \text{Quant}(\Delta_{k+1,i})$ to the central server ;
 - 10 (the central server) ;
 - 11 Set $H_{k+1} = V_k + (np)^{-1} \sum_{i=1}^n B_{k+1,i} \text{Quant}(\Delta_{k+1,i})$;
 - 12 Set $\widehat{S}_{k+1} = \widehat{S}_k + \gamma_{k+1} H_{k+1}$;
 - 13 Set $V_{k+1} = V_k + \alpha n^{-1} \sum_{i=1}^n B_{k+1,i} \text{Quant}(\Delta_{k+1,i})$;
 - 14 Send \widehat{S}_{k+1} and $\mathsf{T}(\widehat{S}_{k+1})$ to the n workers
-

8.2 Notations

Let us introduce a new sequence of filtrations. For any $i \in [n]^*$, we set

$$\mathcal{F}_{0,i} = \mathcal{F}_{0,i}^+ \stackrel{\text{def}}{=} \sigma(\widehat{S}_0; V_{0,i}) \quad \text{and} \quad \mathcal{F}_0 \stackrel{\text{def}}{=} \bigvee_{i=1}^n \mathcal{F}_{0,i} .$$

Then, for all $k \geq 0$,

- (i) $\mathcal{F}_{k+1/3,i} \stackrel{\text{def}}{=} \mathcal{F}_{k,i}^+ \vee \sigma(S_{k+1,i})$,
- (ii) $\mathcal{F}_{k+2/3,i} \stackrel{\text{def}}{=} \mathcal{F}_{k+1/3,i} \vee \sigma(\text{Quant}(\Delta_{k+1,i}))$,
- (iii) $\mathcal{F}_{k+1,i} \stackrel{\text{def}}{=} \mathcal{F}_{k+2/3,i} \vee \sigma(B_{k+1,i})$,
- (iv) $\mathcal{F}_{k+1} \stackrel{\text{def}}{=} \bigvee_{i=1}^n \mathcal{F}_{k+1,i}$,
- (v) $\mathcal{F}_{k+1,i}^+ \stackrel{\text{def}}{=} \mathcal{F}_{k+1,i} \vee \mathcal{F}_{k+1}$.

Note that, with these notations, for $k \geq 0$ and $i \in [n]^*$, the random variables of the FedEM sequence defined in [algorithm 4](#) belong to the filtrations defined above as follows:

- (i) $\widehat{S}_k \in \mathcal{F}_{k,i}^+, \widehat{S}_k \in \mathcal{F}_k,$
- (ii) $S_{k+1,i}, \Delta_{k+1,i} \in \mathcal{F}_{k+1/3,i},$
- (iii) $V_{k+1,i} \in \mathcal{F}_{k+1,i},$
- (iv) $\widehat{S}_{k+1}, H_{k+1}, V_{k+1} \in \mathcal{F}_{k+1}.$

Note also that we have the following inclusions for filtrations: $\mathcal{F}_k \subset \mathcal{F}_{k,i}^+ \subset \mathcal{F}_{k+1/3,i} \subset \mathcal{F}_{k+2/3,i} \subset \mathcal{F}_{k+1,i} \subset \mathcal{F}_{k+1}$ for all $i \in [n]^*$.

8.3 Preliminary results

In this section, we extend Proposition 7, Proposition 8 (that controls the random field H_{k+1}) and Proposition 11 (that controls the memory term $V_{k,i}$). We start by verifying the simple following proposition, that ensures that the global variable V_k corresponds to the mean of the local control variables $(V_{k,i})_{i \in [n]^*}$.

Proposition 12. *For any $k \in [k_{\max}]$,*

$$V_k = \frac{1}{n} \sum_{i=1}^n V_{k,i}.$$

Proof. By definition of V_0 , the property holds true when $k = 0$. Assume this holds true for $k \in [k_{\max} - 1]$. We write

$$\begin{aligned} V_{k+1} &= V_k + \frac{\alpha}{n} \sum_{i=1}^n B_{k+1,i} \text{Quant}(\Delta_{k+1,i}) \\ &= \frac{1}{n} \sum_{i=1}^n V_{k,i} + \frac{1}{n} \sum_{i=1}^n (V_{k+1,i} - V_{k,i}) \\ &= \frac{1}{n} \sum_{i=1}^n V_{k+1,i}. \end{aligned}$$

This concludes the induction. □

We now prove that the unbiased character of H_k is preserved, and we provide a new control on its second order moment. Proposition 13 is Proposition 8 with ω replaced with ω_p . When $p = 1$, Proposition 13 and Proposition 8 are the same.

Proposition 13. *Assume A6, A7 and A8. Set $\sigma^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n \sigma_i^2$. For any $k \in [k_{\max} - 1]$, we have*

$$\mathbb{E}[H_{k+1} | \mathcal{F}_k] = \mathfrak{h}(\widehat{S}_k),$$

and

$$\mathbb{E}[\|H_{k+1} - \mathbb{E}[H_{k+1} | \mathcal{F}_k]\|^2 | \mathcal{F}_k] \leq \frac{\omega_p}{n} \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k] + \frac{\sigma^2}{n},$$

where

$$\omega_p \stackrel{\text{def}}{=} \frac{1-p}{p} (1 + \omega) + \omega. \quad (33)$$

Proof. Let $k \in [k_{\max} - 1]$. By definition, we have

$$H_{k+1} = V_k + \frac{1}{np} \sum_{i=1}^n B_{k+1,i} \text{Quant}(\Delta_{k+1,i})$$

where the Bernoulli random variables $\{B_{k+1,i}, i \in [n]^*\}$ are independent with the same success probability p . By definition of the filtrations, we have $B_{k+1,i} \in \mathcal{F}_{k+1,i}$, $\text{Quant}(\Delta_{k+1,i}) \in \mathcal{F}_{k+2/3,i}$,

$V_k \in \mathcal{F}_k$ and $\Delta_{k+1,i} \in \mathcal{F}_{k+1/3,i}$; and the inclusions $\mathcal{F}_k \subset \mathcal{F}_{k+1/3,i} \subset \mathcal{F}_{k+2/3,i} \subset \mathcal{F}_{k+1,i}$. Therefore,

$$\begin{aligned} \mathbb{E} [H_{k+1} | \mathcal{F}_k] &= V_k + \frac{1}{np} \sum_{i=1}^n \mathbb{E} [\mathbb{E} [B_{k+1,i} | \mathcal{F}_{k+2/3,i}] \text{Quant}(\Delta_{k+1,i}) | \mathcal{F}_k] \\ &= V_k + \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\mathbb{E} [\text{Quant}(\Delta_{k+1,i}) | \mathcal{F}_{k+1/3,i}] | \mathcal{F}_k] = V_k + \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\Delta_{k+1,i} | \mathcal{F}_k] \\ &= V_k + \frac{1}{n} \sum_{i=1}^n \left(\mathbb{E} [S_{k+1,i} | \mathcal{F}_k] - \widehat{S}_k - V_{k,i} \right) \\ &= \frac{1}{n} \sum_{i=1}^n h_i(\widehat{S}_k) = h(\widehat{S}_k), \end{aligned}$$

where we used $\mathbb{E} [B_{k+1,i} | \mathcal{F}_{k+2/3,i}] = p$ (see A8), A6, A7 and Proposition 12. This concludes the proof of the first statement of Proposition 13. For the second point, we write

$$\begin{aligned} H_{k+1} - h(\widehat{S}_k) &= \frac{1}{n} \sum_{i=1}^n \Xi_{k+1,i} \\ \Xi_{k+1,i} &\stackrel{\text{def}}{=} S_{k+1,i} - \mathbb{E} [S_{k+1,i} | \mathcal{F}_{k,i}^+] \\ &\quad + \text{Quant}(\Delta_{k+1,i}) - \mathbb{E} [\text{Quant}(\Delta_{k+1,i}) | \mathcal{F}_{k+1/3,i}] \\ &\quad + \frac{1}{p} (B_{k+1,i} - \mathbb{E} [B_{k+1,i} | \mathcal{F}_{k+2/3,i}]) \text{Quant}(\Delta_{k+1,i}); \end{aligned}$$

note indeed that $h_i(\widehat{S}_k) = \mathbb{E} [S_{k+1,i} | \mathcal{F}_{k,i}^+] - \widehat{S}_k$, $\mathbb{E} [\text{Quant}(\Delta_{k+1,i}) | \mathcal{F}_{k+1/3,i}] = \Delta_{k+1,i}$, $\Delta_{k+1,i} = V_{k,i} + S_{k+1,i} - \widehat{S}_k$, $V_k = n^{-1} \sum_{i=1}^n V_{k,i}$ and $p = \mathbb{E} [B_{k+1,i} | \mathcal{F}_{k+2/3,i}]$. Write $H_{k+1} - h(\widehat{S}_k) = \frac{1}{n} \sum_{i=1}^n \Xi_{k+1,i}$. Since the workers are independent, we have

$$\mathbb{E} [\|H_{k+1} - h(\widehat{S}_k)\|^2 | \mathcal{F}_k] = \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} [\|\Xi_{k+1,i}\|^2 | \mathcal{F}_k].$$

Fix $i \in [n]^*$. $\Xi_{k+1,i}$ is the sum of three terms $\sum_{\ell=1}^3 \Xi_{k+1,i,\ell}$ and observe that for any $\ell \neq \ell'$ we have

$$\mathbb{E} [\langle \Xi_{k+1,i,\ell}, \Xi_{k+1,i,\ell'} \rangle | \mathcal{F}_k] = 0.$$

Therefore $\mathbb{E} [\|\Xi_{k+1,i}\|^2 | \mathcal{F}_k] = \sum_{\ell=1}^3 \mathbb{E} [\|\Xi_{k+1,i,\ell}\|^2 | \mathcal{F}_k]$. We have by A7

$$\mathbb{E} [\|S_{k+1,i} - \mathbb{E} [S_{k+1,i} | \mathcal{F}_{k,i}^+]\|^2 | \mathcal{F}_k] \leq \sigma_i^2;$$

by A6,

$$\mathbb{E} [\|\text{Quant}(\Delta_{k+1,i}) - \mathbb{E} [\text{Quant}(\Delta_{k+1,i}) | \mathcal{F}_{k+1/3,i}]\|^2 | \mathcal{F}_k] \leq \omega \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k];$$

and by A6 and A8

$$\begin{aligned} &\mathbb{E} \left[\frac{1}{p^2} (B_{k+1,i} - \mathbb{E} [B_{k+1,i} | \mathcal{F}_{k+2/3,i}])^2 \|\text{Quant}(\Delta_{k+1,i})\|^2 | \mathcal{F}_k \right] \\ &\leq \frac{1-p}{p} \mathbb{E} [\|\text{Quant}(\Delta_{k+1,i})\|^2 | \mathcal{F}_k] \\ &\leq \frac{1-p}{p} (1+\omega) \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k]. \end{aligned}$$

This concludes the proof. \square

Proposition 14. Assume A7 and set $\sigma^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n \sigma_i^2$. For any $k \in [k_{\max} - 1]$,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k] \leq \frac{1}{n} \sum_{i=1}^n \|V_{k,i} - h_i(\widehat{S}_k)\|^2 + \sigma^2.$$

The proof is on the same lines as the proof of Proposition 10 and is omitted.

Proposition 15 extends Proposition 11: the result is similar but with α replaced with αp and ω by ω_p .

Proposition 15. *Assume A5, A6, A7 and A8; set $L^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n L_i^2$ and $\sigma^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n \sigma_i^2$. Choose $\alpha \in (0, 1/(1 + \omega))$. For any $k \geq 0$, define*

$$G_k \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2.$$

We have, for any $k \in [k_{\max} - 1]$

$$\begin{aligned} \mathbb{E}[G_{k+1} | \mathcal{F}_k] &\leq \left(1 - \frac{\alpha p}{2} + 2\gamma_{k+1}^2 \frac{L^2 \omega_p}{\alpha p n}\right) G_k + 2\gamma_{k+1}^2 \frac{L^2}{\alpha p} \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\quad + 2 \left(\alpha p + \gamma_{k+1}^2 \frac{L^2 \omega_p}{\alpha p n}\right) \sigma^2, \end{aligned}$$

where ω_p is defined in Proposition 13.

Proof. Let $i \in [n]^*$. We follow the same line of the proof as Proposition 11: for any $\beta > 0$, using that $\|a + b\|^2 \leq (1 + \beta^2)\|a\|^2 + (1 + \beta^{-2})\|b\|^2$, we have

$$\begin{aligned} &\mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_{k+1})\|^2 \middle| \mathcal{F}_k \right] \\ &\leq (1 + \beta^{-2}) \mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 \middle| \mathcal{F}_k \right] + (1 + \beta^2) \mathbb{E} \left[\|\mathbf{h}_i(\widehat{S}_k) - \mathbf{h}_i(\widehat{S}_{k+1})\|^2 \middle| \mathcal{F}_k \right] \\ &\stackrel{A5}{\leq} (1 + \beta^{-2}) \mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 \middle| \mathcal{F}_k \right] + (1 + \beta^2) L_i^2 \gamma_{k+1}^2 \mathbb{E} \left[\|H_{k+1}\|^2 \middle| \mathcal{F}_k \right]. \end{aligned}$$

We then provide a control for $\mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 \middle| \mathcal{F}_k \right]$. Recall that:

$$V_{k+1,i} = V_{k,i} + \alpha B_{k+1,i} \text{Quant}(\Delta_{k+1}; i).$$

We write $f(B_{k+1,i}) = f(1)\mathbb{1}_{B_{k+1,i}=1} + f(0)\mathbb{1}_{B_{k+1,i}=0}$ for any measurable positive function f ; and then use $\mathbb{E} \left[\mathbb{1}_{B_{k+1,i}} \middle| \mathcal{F}_{k+2/3,i} \right] = p$ (see A8), $\text{Quant}(\Delta_{k+1}, i)$, $\widehat{S}_k, V_{k,i} \in \mathcal{F}_{k+2/3,i}$. We get

$$\begin{aligned} &\mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 \middle| \mathcal{F}_k \right] \\ &= p \mathbb{E} \left[\|V_{k,i} - \mathbf{h}_i(\widehat{S}_k) - \alpha \text{Quant}(\Delta_{k+1}, i)\|^2 \middle| \mathcal{F}_k \right] + (1 - p) \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 \\ &\stackrel{(21)}{=} p(1 - \alpha) \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 + \alpha p \mathbb{E} \left[\|\mathbf{S}_{k+1,i} - \widehat{S}_k - \mathbf{h}_i(\widehat{S}_k)\|^2 \middle| \mathcal{F}_k \right] \\ &\quad + \alpha p (\alpha(1 + \omega) - 1) \mathbb{E} \left[\|\Delta_{k+1,i}\|^2 \middle| \mathcal{F}_k \right] + (1 - p) \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 \\ &= (1 - \alpha p) \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 \\ &\quad + \alpha p \mathbb{E} \left[\|\mathbf{S}_{k+1,i} - \widehat{S}_k - \mathbf{h}_i(\widehat{S}_k)\|^2 \middle| \mathcal{F}_k \right] + \alpha p (\alpha(1 + \omega) - 1) \mathbb{E} \left[\|\Delta_{k+1,i}\|^2 \middle| \mathcal{F}_k \right]. \end{aligned}$$

The end of the proof is identical to the proof of Proposition 11: we choose $\beta_p > 0$ such that $\beta_p^{-2} = 1$ if $\alpha p \geq 2/3$ and $\beta_p^{-2} = \frac{\alpha p}{2(1 - \alpha p)}$ if $\alpha p \leq 2/3$. We have

$$(1 - \alpha p)(1 + \beta_p^{-2}) \leq 1 - \frac{\alpha p}{2}, \quad (1 + \beta_p^2) \leq \frac{2}{\alpha p}, \quad 1 \leq 1 + \beta_p^{-2} \leq 2;$$

and this yields

$$\begin{aligned} &\mathbb{E} \left[\|V_{k+1,i} - \mathbf{h}_i(\widehat{S}_{k+1})\|^2 \middle| \mathcal{F}_k \right] \leq \left(1 - \frac{\alpha p}{2}\right) \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 \\ &\quad + 2\alpha p \mathbb{E} \left[\|\mathbf{S}_{k+1,i} - \bar{\mathbf{s}}_i \circ \mathbf{T}(\widehat{S}_k)\|^2 \middle| \mathcal{F}_k \right] + \alpha p (\alpha(1 + \omega) - 1) \mathbb{E} \left[\|\Delta_{k+1,i}\|^2 \middle| \mathcal{F}_k \right] \\ &\quad + \frac{2}{\alpha p} L_i^2 \gamma_{k+1}^2 \mathbb{E} \left[\|H_{k+1}\|^2 \middle| \mathcal{F}_k \right]. \end{aligned}$$

By definition of the conditional expectation and Proposition 13 we have

$$\begin{aligned}\mathbb{E} [\|H_{k+1}\|^2 | \mathcal{F}_k] &= \|\mathbb{E} [H_{k+1} | \mathcal{F}_k]\|^2 + \mathbb{E} [\|H_{k+1} - \mathbb{E} [H_{k+1} | \mathcal{F}_k]\|^2 | \mathcal{F}_k] \\ &= \|\mathbf{h}(\widehat{S}_k)\|^2 + \mathbb{E} [\|H_{k+1} - \mathbf{h}(\widehat{S}_k)\|^2 | \mathcal{F}_k] .\end{aligned}$$

Since $(\alpha(1 + \omega) - 1) \leq 0$, using A7 and Proposition 13 again, we get:

$$\mathbb{E} [G_{k+1} | \mathcal{F}_k] \leq \left(1 - \frac{\alpha p}{2}\right) G_k + 2\alpha p \sigma^2 + \frac{2}{\alpha p} L^2 \gamma_{k+1}^2 \frac{1}{n} \left(\sigma^2 + \omega_p \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k] \right) .$$

Finally, from Proposition 14,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\Delta_{k+1,i}\|^2 | \mathcal{F}_k] \leq G_k + \sigma^2 .$$

This concludes the proof. □

8.4 Proof of Theorem 3

Throughout this proof, set

$$\omega_p \stackrel{\text{def}}{=} \frac{1-p}{p}(1+\omega) + \omega .$$

Step 1: Upper bound on the decrement. Let $k \geq 0$. Following the same lines as in the proof of Theorem 1, we have

$$\begin{aligned}\mathbb{E} [\mathbf{W}(\widehat{S}_{k+1}) | \mathcal{F}_k] &\leq \mathbf{W}(\widehat{S}_k) - \gamma_{k+1} v_{\min} \left(1 - \gamma_{k+1} \frac{L\dot{\mathbf{W}}}{2v_{\min}}\right) \|\mathbf{h}(\widehat{S}_k)\|^2 + \gamma_{k+1}^2 \frac{L\dot{\mathbf{W}}}{2} \mathbb{E} [\|H_{k+1} - \mathbb{E} [H_{k+1} | \mathcal{F}_k]\|^2 | \mathcal{F}_k] .\end{aligned}$$

Applying Proposition 13 and Proposition 14, we obtain that

$$\begin{aligned}\mathbb{E} [\mathbf{W}(\widehat{S}_{k+1}) | \mathcal{F}_k] &\leq \mathbf{W}(\widehat{S}_k) - \gamma_{k+1} v_{\min} \left(1 - \gamma_{k+1} \frac{L\dot{\mathbf{W}}}{2v_{\min}}\right) \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\quad + \gamma_{k+1}^2 \frac{L\dot{\mathbf{W}}}{2} \frac{\omega_p}{n} G_k + \gamma_{k+1}^2 \frac{L\dot{\mathbf{W}}}{2n} (1 + \omega_p) \sigma^2 , \quad (34)\end{aligned}$$

where

$$G_k \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|V_{k,i} - \mathbf{h}_i(\widehat{S}_k)\|^2 .$$

Step 2: Maximal learning rate γ_{k+1} when $\omega \neq 0$. From Proposition 11, for any non-increasing positive sequence $\{\gamma_k, k \in [k_{\max} - 1]\}$ such that

$$\gamma_{k+1}^2 \leq \frac{\alpha^2 p^2}{8L^2} \frac{n}{\omega_p} ,$$

and for any positive sequence $\{C_k, k \in [k_{\max} - 1]\}$, it holds

$$\begin{aligned}C_{k+1} \mathbb{E} [G_{k+1} | \mathcal{F}_k] &\leq C_{k+1} \left(1 - \frac{\alpha p}{4}\right) G_k \\ &\quad + C_{k+1} \gamma_{k+1}^2 \frac{2}{\alpha p} L^2 \|\mathbf{h}(\widehat{S}_k)\|^2 + 2C_{k+1} \left(\alpha p + \gamma_{k+1}^2 \frac{L^2}{\alpha p} \frac{1 + \omega_p}{n}\right) \sigma^2 . \quad (35)\end{aligned}$$

Combining equations (34) and (35), we thus have

$$\begin{aligned} \mathbb{E}[\mathbb{W}(\widehat{S}_{k+1})|\mathcal{F}_k] + C_{k+1}\mathbb{E}[G_{k+1}|\mathcal{F}_k] &\leq \mathbb{W}(\widehat{S}_k) + C_k G_k \\ &\quad - \gamma_{k+1}v_{\min} \left(1 - \gamma_{k+1} \frac{L\dot{\mathbb{W}}}{2v_{\min}} - \frac{C_{k+1}}{v_{\min}} \gamma_{k+1} \frac{2}{\alpha p} L^2 \right) \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\quad + \left(\gamma_{k+1}^2 \frac{L\dot{\mathbb{W}}}{2} \frac{\omega_p}{n} - C_k + C_{k+1} - C_{k+1} \frac{\alpha p}{4} \right) G_k \\ &\quad + \left\{ 2\alpha p C_{k+1} + \gamma_{k+1}^2 \frac{(1 + \omega_p)}{n} \left(\frac{L\dot{\mathbb{W}}}{2} + 2C_{k+1} \frac{L^2}{\alpha p} \right) \right\} \sigma^2. \end{aligned}$$

We choose the sequence $\{C_k\}$ as follows:

$$C_k \stackrel{\text{def}}{=} \gamma_k^2 \frac{2L\dot{\mathbb{W}}}{\alpha p} \frac{\omega_p}{n};$$

the sequence satisfies $C_{k+1} \leq C_k$ (since $\gamma_{k+1} \leq \gamma_k$) and $\gamma_{k+1}^2 L\dot{\mathbb{W}}\omega_p/(2n) \leq C_{k+1}\alpha p/4$. By convention, $\gamma_0 \in [\gamma_1, +\infty)$. Therefore

$$\begin{aligned} \mathbb{E}[\mathbb{W}(\widehat{S}_{k+1})|\mathcal{F}_k] + \gamma_{k+1}^2 \frac{2L\dot{\mathbb{W}}}{\alpha p} \frac{\omega_p}{n} \mathbb{E}[G_{k+1}|\mathcal{F}_k] &\leq \mathbb{W}(\widehat{S}_k) + \gamma_k^2 \frac{2L\dot{\mathbb{W}}}{\alpha p} \frac{\omega_p}{n} G_k \\ &\quad - \gamma_{k+1}v_{\min} \left(1 - \gamma_{k+1} \frac{L\dot{\mathbb{W}}}{2v_{\min}} \left\{ 1 + 8\gamma_{k+1}^2 \frac{\omega_p}{\alpha^2 p^2 n} L^2 \right\} \right) \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\quad + 4\gamma_{k+1}^2 L\dot{\mathbb{W}} \frac{\omega_p}{n} \left\{ 1 + \frac{(1 + \omega_p)}{8\omega_p} \left(1 + \gamma_{k+1}^2 8 \frac{L^2}{\alpha^2 p^2} \frac{\omega_p}{n} \right) \right\} \sigma^2. \end{aligned}$$

Step 3: Computing the expectation. Let us apply the expectations, sum from $k = 0$ to $k = k_{\max} - 1$, and divide by k_{\max} . This yields

$$\begin{aligned} &\frac{v_{\min}}{k_{\max}} \sum_{k=0}^{k_{\max}-1} \gamma_{k+1} \left(1 - \gamma_{k+1} \frac{L\dot{\mathbb{W}}}{2v_{\min}} \left\{ 1 + 8\gamma_{k+1}^2 \frac{\omega_p}{\alpha^2 p^2 n} L^2 \right\} \right) \|\mathbf{h}(\widehat{S}_k)\|^2 \\ &\leq k_{\max}^{-1} \left\{ \mathbb{W}(\widehat{S}_0) + \gamma_0^2 \frac{2L\dot{\mathbb{W}}}{\alpha} \frac{\omega_p}{n} G_0 - \mathbb{E}[\mathbb{W}(\widehat{S}_{k_{\max}})] - \gamma_{k_{\max}}^2 \frac{2L\dot{\mathbb{W}}}{\alpha p} \frac{\omega_p}{n} \mathbb{E}[G_{k_{\max}}] \right\} \\ &\quad + 4L\dot{\mathbb{W}} \frac{\omega_p}{n} \frac{1}{k_{\max}} \sum_{k=0}^{k_{\max}-1} \gamma_{k+1}^2 \left\{ 1 + \frac{(1 + \omega_p)}{8\omega_p} \left(1 + \gamma_{k+1}^2 8 \frac{L^2}{\alpha^2 p^2} \frac{\omega_p}{n} \right) \right\} \sigma^2. \end{aligned}$$

We now focus on the case when $\gamma_{k+1} = \gamma$ for any $k \geq 0$. Denote by K a uniform random variable on $[k_{\max} - 1]$, independent of the path $\{\widehat{S}_k, k \in [k_{\max}]\}$. Since $\gamma^2 \leq \alpha^2 p^2 n / (8L^2 \omega_p)$, we have

$$1 + 8\gamma^2 \frac{\omega_p}{\alpha^2 p^2 n} L^2 \leq 2.$$

This yields

$$\begin{aligned} &v_{\min} \gamma \left(1 - \gamma \frac{L\dot{\mathbb{W}}}{v_{\min}} \right) \mathbb{E}[\|\mathbf{h}(\widehat{S}_K)\|^2] \\ &\leq k_{\max}^{-1} \left\{ \mathbb{W}(\widehat{S}_0) + \gamma^2 \frac{2L\dot{\mathbb{W}}}{\alpha p} \frac{\omega_p}{n} G_0 - \mathbb{E}[\mathbb{W}(\widehat{S}_{k_{\max}})] - \gamma^2 \frac{2L\dot{\mathbb{W}}}{\alpha p} \frac{\omega_p}{n} \mathbb{E}[G_{k_{\max}}] \right\} \\ &\quad + 4L\dot{\mathbb{W}} \frac{\omega_p}{n} \gamma^2 \left\{ 1 + \frac{(1 + \omega_p)}{4\omega_p} \right\} \sigma^2. \end{aligned}$$

Note that $4(1 + (1 + \omega_p)/(4\omega_p)) = (5\omega_p + 1)/\omega_p$.

Step 4. Conclusion (when $\omega \neq 0$) By choosing $V_{0,i} = \mathbf{h}_i$ for any $i \in [n]^*$, we have $G_0 = 0$. The roots of $\gamma \mapsto \gamma(1 - \gamma L\dot{\mathbb{W}}/v_{\min})$ are 0 and $v_{\min}/L\dot{\mathbb{W}}$ and its maximum is reached at $v_{\min}/(2L\dot{\mathbb{W}})$: this function is increasing on $(0, v_{\min}/(2L\dot{\mathbb{W}})]$. We therefore choose $\gamma \in (0, \gamma_{\max}(\alpha)]$ where

$$\gamma_{\max}(\alpha) \stackrel{\text{def}}{=} \min \left(\frac{v_{\min}}{2L\dot{\mathbb{W}}}; \frac{\alpha p}{2\sqrt{2}L} \frac{\sqrt{n}}{\sqrt{\omega_p}} \right)$$

Finally, since $\alpha \in (0, 1/(1 + \omega)]$, we choose $\alpha = 1/(1 + \omega)$. This yields

$$\gamma_{\max} \stackrel{\text{def}}{=} \min \left(\frac{v_{\min}}{2L_{\text{W}}}; \frac{p}{2\sqrt{2}L} \frac{\sqrt{n}}{\sqrt{\omega_p}(1 + \omega)} \right).$$

9 Convergence Analysis of VR-FedEM

The assumptions A1 to A3 hold throughout this section. We will use the notations

$$L_i^2 \stackrel{\text{def}}{=} m^{-1} \sum_{j=1}^m L_{ij}^2, \quad L^2 \stackrel{\text{def}}{=} n^{-1} \sum_{i=1}^n L_i^2, \quad (36)$$

where L_{ij} is defined in A9, and

$$h_i(s) \stackrel{\text{def}}{=} \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(s) - s, \quad h(s) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n h_i(s).$$

9.1 Notations and elementary result

Let us define the following filtrations: for any $i \in [n]^*$ and $t \in [k_{\text{out}}]^*$, $k \in [k_{\text{max}} - 1]$, set

$$\begin{aligned} \mathcal{F}_{1,0,i} &= \mathcal{F}_{1,0,i}^+ \stackrel{\text{def}}{=} \sigma \left(\widehat{S}_{\text{init}}; V_{1,0,i} \right), & \mathcal{F}_{1,0} &\stackrel{\text{def}}{=} \bigvee_{i=1}^n \mathcal{F}_{1,0,i}, \\ \mathcal{F}_{t,k+1/2,i} &\stackrel{\text{def}}{=} \mathcal{F}_{t,k,i}^+ \vee \sigma \left(\mathcal{B}_{t,k+1,i} \right), & \mathcal{F}_{t,k+1,i} &\stackrel{\text{def}}{=} \mathcal{F}_{t,k+1/2,i} \vee \sigma \left(\text{Quant}(\Delta_{t,k+1,i}) \right), \\ \mathcal{F}_{t,k+1} &\stackrel{\text{def}}{=} \bigvee_{i=1}^n \mathcal{F}_{t,k+1,i}, & \mathcal{F}_{t,k+1,i}^+ &\stackrel{\text{def}}{=} \mathcal{F}_{t,k+1}. \end{aligned}$$

With these notations, for $t \in [k_{\text{out}}]^*$, $k \in [k_{\text{max}} - 1]$ and $i \in [n]^*$, $\widehat{S}_{t,k+1} \in \mathcal{F}_{t,k+1,i}^+$, $S_{t,k+1,i} \in \mathcal{F}_{t,k+1/2,i}$, $\Delta_{t,k+1,i} \in \mathcal{F}_{t,k+1/2,i}$, $V_{t,k+1,i} \in \mathcal{F}_{t,k+1,i}$, $\widehat{S}_{t,k+1} \in \mathcal{F}_{t,k+1}$, $H_{t,k+1} \in \mathcal{F}_{t,k+1}$, and $V_{t,k+1} \in \mathcal{F}_{t,k+1}$.

9.2 Preliminary results

9.2.1 Results on the minibatch $\mathcal{B}_{t,k+1}$

The proof of the following proposition is given in [10, Lemma 4]. It establishes the bias and the variance of the sum along the random set of indices $\mathcal{B}_{t,k+1}$ conditionally to the past.

Proposition 16. *Let \mathcal{B} be a minibatch of size b , sampled at random (with or without replacement) from $[m]^*$. It holds for any $i \in [n]^*$ and $s \in \mathbb{R}^q$,*

$$\mathbb{E} \left[\frac{1}{b} \sum_{j \in \mathcal{B}} \bar{s}_{ij} \circ \mathbb{T}(s) \right] = \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(s);$$

and for any $s, s' \in \mathbb{R}^q$,

$$\begin{aligned} \mathbb{E} \left[\left\| \frac{1}{b} \sum_{j \in \mathcal{B}} \{ \bar{s}_{ij} \circ \mathbb{T}(s) - s \} - \{ \bar{s}_{ij} \circ \mathbb{T}(s') - s' \} \right\|^2 \right] \\ - \frac{1}{m} \sum_{j=1}^m \left\| \{ \bar{s}_{ij} \circ \mathbb{T}(s) - s \} - \{ \bar{s}_{ij} \circ \mathbb{T}(s') - s' \} \right\|^2 \leq \frac{L_i^2}{b} \|s - s'\|^2. \end{aligned}$$

9.2.2 Results on the statistics $S_{t,k,i}$

Proposition 17 shows that for $k \geq 1$, $S_{t,k+1,i}$ is a biased approximation of $m^{-1} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k})$; and this bias is canceled at the beginning of each outer loop since $S_{t,1,i} = m^{-1} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,0})$. Corollary 18 establishes an upper bound for the conditional variance and the mean squared error of $S_{t,k+1,i}$.

Let us comment the definition of $S_{t,k+1,i}$. For any $t \in [k_{\text{out}}]^*$, $k \in [k_{\text{in}} - 1]$ and $i \in [n]^*$,

$$S_{t,k+1,i} = \frac{1}{b} \sum_{j \in \mathcal{B}_{t,k+1,i}} \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) + \Upsilon_{t,k+1,i}, \quad \Upsilon_{t,k+1,i} \stackrel{\text{def}}{=} S_{t,k,i} - \frac{1}{b} \sum_{j \in \mathcal{B}_{t,k+1,i}} \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k-1}).$$

It is easily seen that

$$\Upsilon_{t,k+1,i} = \Upsilon_{t,k,i} + \frac{1}{b} \sum_{j \in \mathcal{B}_{t,k,i}} \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k-1}) - \frac{1}{b} \sum_{j \in \mathcal{B}_{t,k+1,i}} \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k-1}),$$

and since $\Upsilon_{t,1,i} = S_{t,0,i} - b^{-1} \sum_{j \in \mathcal{B}_{t,1,i}} \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,-1})$, we have by using Proposition 17,

$$\begin{aligned} \Upsilon_{t,k,i} &= \sum_{\ell=1}^k \left\{ \frac{1}{b} \sum_{j \in \mathcal{B}_{t,\ell,i}} \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,\ell-1}) - \frac{1}{b} \sum_{j \in \mathcal{B}_{t,\ell+1,i}} \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,\ell-1}) \right\} \\ &\quad + \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,-1}) - \frac{1}{b} \sum_{j \in \mathcal{B}_{t,1,i}} \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,-1}). \end{aligned}$$

We have $\mathbb{E}[\Upsilon_{t,k,i} | \mathcal{F}_{t,0}] = 0$ but conditionally to the past $\mathcal{F}_{t,k-1,i}^+$, the variable $\Upsilon_{t,k,i}$ is *not* centered.

Proposition 17. For any $t \in [k_{\text{out}}]^*$ and $i \in [n]^*$,

$$S_{t,1,i} - \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,0}) = S_{t,0,i} - \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,-1}) = 0.$$

For any $t \in [k_{\text{out}}]^*$, $k \in [k_{\text{in}} - 1]$ and $i \in [n]^*$, we have

$$\mathbb{E} \left[S_{t,k+1,i} \middle| \mathcal{F}_{t,k,i}^+ \right] - \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) = S_{t,k,i} - \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k-1}).$$

Proof. Let $t \in [k_{\text{out}}]^*$ and $i \in [n]^*$. We have by definition of $S_{t,1,i}$ and $S_{t,0,i}$

$$S_{t,1,i} = S_{t,0,i} + b^{-1} \sum_{j \in \mathcal{B}_{t,1,i}} \left(\bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,0}) - \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,-1}) \right) = S_{t,0,i} = \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,0})$$

where we used that $\widehat{S}_{t,0} = \widehat{S}_{t,-1}$.

Let $k \in [k_{\text{in}} - 1]$. By definition of $S_{t,k+1,i}$, we have

$$S_{t,k+1,i} - S_{t,k,i} = b^{-1} \sum_{j \in \mathcal{B}_{t,k+1,i}} \left(\bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) - \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k-1}) \right).$$

Since $\widehat{S}_{t,k}, \widehat{S}_{t,k-1} \in \mathcal{F}_{t,k,i}^+$, we have by Proposition 16

$$\begin{aligned} \mathbb{E} \left[b^{-1} \sum_{j \in \mathcal{B}_{t,k+1,i}} \left(\bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) - \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k-1}) \right) \middle| \mathcal{F}_{t,k,i}^+ \right] \\ = \frac{1}{m} \sum_{j=1}^m \left(\bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) - \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k-1}) \right) \end{aligned}$$

and the proof follows. \square

Corollary 18 (of Proposition 17). *Assume A9. For any $t \in [k_{\text{out}}]^*$, $k \in [k_{\text{in}} - 1]$ and $i \in [n]^*$,*

$$\begin{aligned} \mathbb{E} \left[\left\| \mathbf{S}_{t,k+1,i} - \mathbb{E} \left[\mathbf{S}_{t,k+1,i} \mid \mathcal{F}_{t,k,i} \right] \right\|^2 \mid \mathcal{F}_{t,k} \right] &\leq \frac{L_i^2}{\mathbf{b}} \gamma_{t,k}^2 \|H_{t,k}\|^2, \\ \mathbb{E} \left[\left\| \mathbf{S}_{t,k+1,i} - \frac{1}{m} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k}) \right\|^2 \mid \mathcal{F}_{t,0} \right] &\leq \frac{L_i^2}{\mathbf{b}} \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} \left[\|H_{t,\ell}\|^2 \mid \mathcal{F}_{t,0} \right]. \end{aligned}$$

By convention, $H_{t,0} = 0$ and $\sum_{\ell=1}^0 a_\ell = 0$.

Proof. Note that $\widehat{\mathbf{S}}_{t,k}, \widehat{\mathbf{S}}_{t,k-1} \in \mathcal{F}_{t,k}$. By Proposition 17, we have

$$\begin{aligned} &\mathbb{E} \left[\left\| \mathbf{S}_{t,k+1,i} - \mathbb{E} \left[\mathbf{S}_{t,k+1,i} \mid \mathcal{F}_{t,k,i}^+ \right] \right\|^2 \mid \mathcal{F}_{t,k} \right] \\ &= \mathbb{E} \left[\left\| \frac{1}{\mathbf{b}} \sum_{j \in \mathcal{B}_{t,k+1,i}} \left(\bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k}) - \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k-1}) \right) - \frac{1}{m} \sum_{j=1}^m \left(\bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k}) - \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k-1}) \right) \right\|^2 \mid \mathcal{F}_{t,k} \right]. \end{aligned}$$

By Proposition 16, it holds

$$\mathbb{E} \left[\left\| \mathbf{S}_{t,k+1,i} - \mathbb{E} \left[\mathbf{S}_{t,k+1,i} \mid \mathcal{F}_{t,k,i} \right] \right\|^2 \mid \mathcal{F}_{t,k} \right] \leq \frac{L_i^2}{\mathbf{b}} \|\widehat{\mathbf{S}}_{t,k} - \widehat{\mathbf{S}}_{t,k-1}\|^2 = \frac{L_i^2}{\mathbf{b}} \gamma_{t,k}^2 \|H_{t,k}\|^2;$$

with the convention that $H_{t,0} = 0$ since $\widehat{\mathbf{S}}_{t,0} = \widehat{\mathbf{S}}_{t,-1}$. The proof of the first statement is concluded.

For the second statement, by definition of the conditional expectation and since $\widehat{\mathbf{S}}_{t,k} \in \mathcal{F}_{t,k} \subset \mathcal{F}_{t,k,i}^+$, it holds

$$\begin{aligned} \mathbb{E} \left[\left\| \mathbf{S}_{t,k+1,i} - \frac{1}{m} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k}) \right\|^2 \mid \mathcal{F}_{t,k} \right] &= \mathbb{E} \left[\left\| \mathbf{S}_{t,k+1,i} - \mathbb{E} \left[\mathbf{S}_{t,k+1,i} \mid \mathcal{F}_{t,k,i}^+ \right] \right\|^2 \mid \mathcal{F}_{t,k} \right] \\ &\quad + \mathbb{E} \left[\left\| \mathbb{E} \left[\mathbf{S}_{t,k+1,i} \mid \mathcal{F}_{t,k,i}^+ \right] - \frac{1}{m} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k}) \right\|^2 \mid \mathcal{F}_{t,k} \right]. \end{aligned}$$

By Proposition 17,

$$\left\| \mathbb{E} \left[\mathbf{S}_{t,k+1,i} \mid \mathcal{F}_{t,k,i}^+ \right] - \frac{1}{m} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k}) \right\|^2 = \left\| \mathbf{S}_{t,k,i} - \frac{1}{m} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k-1}) \right\|^2.$$

Hence, by using $\mathbf{S}_{t,1,i} - m^{-1} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,0}) = 0$ (see Proposition 17), we have

$$\begin{aligned} &\mathbb{E} \left[\left\| \mathbf{S}_{t,k+1,i} - \frac{1}{m} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k}) \right\|^2 \mid \mathcal{F}_{t,0} \right] \\ &\leq \frac{L_i^2}{\mathbf{b}} \gamma_{t,k}^2 \mathbb{E} \left[\|H_{t,k}\|^2 \mid \mathcal{F}_{t,0} \right] + \mathbb{E} \left[\left\| \mathbf{S}_{t,k,i} - \frac{1}{m} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{\mathbf{S}}_{t,k-1}) \right\|^2 \mid \mathcal{F}_{t,0} \right] \\ &\leq \frac{L_i^2}{\mathbf{b}} \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} \left[\|H_{t,\ell}\|^2 \mid \mathcal{F}_{t,0} \right]. \end{aligned}$$

□

9.2.3 Results on $\Delta_{t,k+1,i}$

Proposition 19 provides an upper bound for the mean value of the conditional variance of $\Delta_{t,k+1,i}$, and for its L_2 -moment. Proposition 20 prepares the control of the variance of the random field $H_{t,k+1}$ upon noting that

$$H_{t,k+1} - \mathbb{E} \left[H_{t,k+1} \mid \mathcal{F}_{t,k} \right] = \frac{1}{n} \sum_{i=1}^n \left(\text{Quant}(\Delta_{t,k+1,i}) - \mathbb{E} \left[\Delta_{t,k+1,i} \mid \mathcal{F}_{t,k} \right] \right).$$

Proposition 19. *Assume A9. For any $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}} - 1]$,*

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,0}] \\ & \leq 2 \frac{L^2}{\mathfrak{b}} \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} [\|H_{t,\ell}\|^2 | \mathcal{F}_{t,0}] + \frac{2}{n} \sum_{i=1}^n \mathbb{E} [\|\mathbf{h}_i(\widehat{S}_{t,k}) - V_{t,k,i}\|^2 | \mathcal{F}_{t,0}]. \end{aligned}$$

In addition,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\Delta_{t,k+1,i} - \mathbb{E}[\Delta_{t,k+1,i} | \mathcal{F}_{t,k}]\|^2 | \mathcal{F}_{t,k}] \leq \frac{L^2}{\mathfrak{b}} \gamma_{t,k}^2 \|H_{t,k}\|^2.$$

Proof. Let $i \in [n]^*$, $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}} - 1]$. We write

$$\Delta_{t,k+1,i} = S_{t,k+1,i} - \frac{1}{m} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) + \mathbf{h}_i(\widehat{S}_{t,k}) - V_{t,k,i}.$$

When $k = 0$, we have $S_{t,1,i} - \frac{1}{m} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{S}_{t,0}) = 0$ (see Proposition 17) so that $\Delta_{t,1,i} = \mathbf{h}_i(\widehat{S}_{t,0}) - V_{t,0,i}$. For $k \geq 1$, we write

$$\begin{aligned} \mathbb{E} [\|\Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,0}] & \leq 2 \mathbb{E} \left[\left\| S_{t,k+1,i} - \frac{1}{m} \sum_{j=1}^m \bar{\mathbf{s}}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) \right\|^2 \middle| \mathcal{F}_{t,0} \right] \\ & \quad + 2 \mathbb{E} [\|\mathbf{h}_i(\widehat{S}_{t,k}) - V_{t,k,i}\|^2 | \mathcal{F}_{t,0}] \end{aligned}$$

and the proof of the first statement is concluded by Corollary 18.

By definition of $\Delta_{t,k+1,i}$, it holds

$$\Delta_{t,k+1,i} - \mathbb{E}[\Delta_{t,k+1,i} | \mathcal{F}_{t,k}] = S_{t,k+1,i} - \mathbb{E}[S_{t,k+1,i} | \mathcal{F}_{t,k}]. \quad (37)$$

The proof is concluded by (37) and Corollary 18. \square

Proposition 20. *Assume A6 and A9. For any $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}} - 1]$,*

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\text{Quant}(\Delta_{t,k+1,i}) - \mathbb{E}[\Delta_{t,k+1,i} | \mathcal{F}_{t,k}]\|^2 | \mathcal{F}_{t,0}] & \leq \frac{\omega}{n} \sum_{i=1}^n \mathbb{E} [\|\Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,0}] \\ & \quad + \frac{L^2}{\mathfrak{b}} \gamma_{t,k}^2 \mathbb{E} [\|H_{t,k}\|^2 | \mathcal{F}_{t,0}]. \end{aligned}$$

Proof. Let $i \in [n]^*$, $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}} - 1]$. We write

$$\text{Quant}(\Delta_{t,k+1,i}) - \mathbb{E}[\Delta_{t,k+1,i} | \mathcal{F}_{t,k}] = \text{Quant}(\Delta_{t,k+1,i}) - \Delta_{t,k+1,i} + \Delta_{t,k+1,i} - \mathbb{E}[\Delta_{t,k+1,i} | \mathcal{F}_{t,k}];$$

and use the property

$$\begin{aligned} \mathbb{E} [\|\text{Quant}(\Delta_{t,k+1,i}) - \mathbb{E}[\Delta_{t,k+1,i} | \mathcal{F}_{t,k}]\|^2 | \mathcal{F}_{t,0}] & = \mathbb{E} [\|\text{Quant}(\Delta_{t,k+1,i}) - \Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,0}] \\ & \quad + \mathbb{E} [\|\Delta_{t,k+1,i} - \mathbb{E}[\Delta_{t,k+1,i} | \mathcal{F}_{t,k}]\|^2 | \mathcal{F}_{t,0}]. \end{aligned}$$

By A6 and $\mathcal{F}_{t,k} \subset \mathcal{F}_{t,k+1/2,i}$, we have

$$\begin{aligned} & \mathbb{E} [\|\text{Quant}(\Delta_{t,k+1,i}) - \Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,0}] \\ & = \mathbb{E} [\mathbb{E} [\|\text{Quant}(\Delta_{t,k+1,i}) - \Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,k+1/2,i}] | \mathcal{F}_{t,0}] \leq \omega \mathbb{E} [\|\Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,0}]; \end{aligned}$$

in addition, by Proposition 19,

$$n^{-1} \sum_{i=1}^n \mathbb{E} [\|\Delta_{t,k+1,i} - \mathbb{E}[\Delta_{t,k+1,i} | \mathcal{F}_{t,k}]\|^2 | \mathcal{F}_{t,0}] \leq \frac{L^2}{\mathfrak{b}} \gamma_{t,k}^2 \mathbb{E} [\|H_{t,k}\|^2 | \mathcal{F}_{t,0}].$$

This concludes the proof. \square

9.2.4 Results on the memory terms $V_{t,k+1,i}$

Lemma 21 proves that the memory term $V_{t,k+1}$ computed by the central server is the mean value of the local $V_{t,k+1,i}$ computed by each worker $\#i$. Proposition 22 establishes a contraction-like inequality on the mean quantity $n^{-1} \sum_{i=1}^n \|V_{t,k+1,i} - h_i(\widehat{S}_{t,k+1})\|^2$ thus providing the intuition that $V_{t,k+1,i}$ approximates $h_i(\widehat{S}_{t,k+1})$.

Lemma 21. *For any $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}} - 1]$,*

$$V_{t,k+1} = \frac{1}{n} \sum_{i=1}^n V_{t,k+1,i}, \quad V_{t,0} = \frac{1}{n} \sum_{i=1}^n V_{t,0,i}.$$

Proof. The proof is by induction on t and k . Consider the case $t = 1$. When $k = 0$, the property holds true by Line 1 in algorithm 2. Assume that the property holds for $k \leq k_{\text{in}} - 2$. Then by definition of $V_{1,k+1}$ and by the induction assumption:

$$\begin{aligned} V_{1,k+1} &= V_{1,k} + \alpha \frac{1}{n} \sum_{i=1}^n \text{Quant}(\Delta_{1,k+1,i}) = \frac{1}{n} \sum_{i=1}^n (V_{1,k,i} + \alpha \text{Quant}(\Delta_{1,k+1,i})) \\ &= \frac{1}{n} \sum_{i=1}^n V_{1,k+1,i}. \end{aligned}$$

By Lines 18 and 21 in algorithm 2 and by the induction on k , we obtain

$$V_{2,0} = V_{1,k_{\text{in}}} = \frac{1}{n} \sum_{i=1}^n V_{1,k_{\text{in}},i} = \frac{1}{n} \sum_{i=1}^n V_{2,0,i}.$$

Assume that for $t \in [k_{\text{out}} - 1]^*$ we have $V_{t,0} = n^{-1} \sum_{i=1}^n V_{t,0,i}$. As in the case $t = 1$, we prove by induction on k that for any $k \in [k_{\text{in}} - 1]$, $V_{t,k+1} = n^{-1} \sum_{i=1}^n V_{t,k+1,i}$ (details are omitted). This implies, by using Lines 18 and 21 of algorithm 2, that

$$V_{t+1,0} = V_{t,k_{\text{in}}} = \frac{1}{n} \sum_{i=1}^n V_{t,k_{\text{in}},i} = \frac{1}{n} \sum_{i=1}^n V_{t+1,0,i}.$$

This concludes the induction. \square

Proposition 22. *Assume A6 and A9. Let $\alpha \in (0, (1 + \omega)^{-1})$. For any $t \in [k_{\text{out}}]^*$, $k \in [k_{\text{in}} - 1]$ and $i \in [n]^*$, it holds*

$$\mathbb{E} [V_{t,k+1,i} | \mathcal{F}_{t,k+1/2,i}] = (1 - \alpha) V_{t,k,i} + \alpha (S_{t,k+1,i} - \widehat{S}_{t,k}),$$

Define for $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}}]$

$$G_{t,k} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|V_{t,k,i} - h_i(\widehat{S}_{t,k})\|^2.$$

We have

$$\begin{aligned} \mathbb{E} [G_{t,k+1} | \mathcal{F}_{t,0}] &\leq (1 - \alpha/2) \mathbb{E} [G_{t,k} | \mathcal{F}_{t,0}] \\ &\quad + \frac{2}{\alpha} L^2 \gamma_{t,k+1}^2 \mathbb{E} [\|H_{t,k+1}\|^2 | \mathcal{F}_{t,0}] + 2\alpha \frac{L^2}{b} \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} [\|H_{t,\ell}\|^2 | \mathcal{F}_{t,0}] \\ &\quad + \alpha (\alpha(1 + \omega) - 1) \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|\Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,0}]. \end{aligned}$$

Proof. Let $t \in [k_{\text{out}}]^*$, $k \in [k_{\text{in}} - 1]$ and $i \in [n]^*$. By definition of $V_{t,k+1,i}$, $\Delta_{t,k+1,i}$ and by A6, it holds

$$\begin{aligned} \mathbb{E} [V_{t,k+1,i} | \mathcal{F}_{t,k+1/2,i}] &= V_{t,k,i} + \alpha \mathbb{E} [\text{Quant}(\Delta_{t,k+1,i}) | \mathcal{F}_{t,k+1/2,i}] \\ &= V_{t,k,i} + \alpha (S_{t,k+1,i} - \widehat{S}_{t,k} - V_{t,k,i}). \end{aligned}$$

This concludes the proof of the first statement. For the second statement, we write for any $\beta > 0$:

$$\begin{aligned} \|V_{t,k+1,i} - \mathbf{h}_i(\widehat{S}_{t,k+1})\|^2 &\leq (1 + \beta^2) \|\mathbf{h}_i(\widehat{S}_{t,k+1}) - \mathbf{h}_i(\widehat{S}_{t,k})\|^2 + (1 + \beta^{-2}) \|V_{t,k+1,i} - \mathbf{h}_i(\widehat{S}_{t,k})\|^2 \\ &\leq (1 + \beta^2) L_i^2 \gamma_{t,k+1}^2 \|H_{t,k+1}\|^2 + (1 + \beta^{-2}) \|V_{t,k+1,i} - \mathbf{h}_i(\widehat{S}_{t,k})\|^2, \end{aligned} \quad (38)$$

where we used A9 and the definition of $\widehat{S}_{t,k+1}$ in the last inequality. For any $s \in \mathbb{R}^q$

$$\begin{aligned} \mathbb{E} [\|V_{t,k+1,i} - s\|^2 | \mathcal{F}_{t,k+1/2,i}] &= \mathbb{E} [\|V_{t,k+1,i} - \mathbb{E}[V_{t,k+1,i} | \mathcal{F}_{t,k+1/2,i}]\|^2 | \mathcal{F}_{t,k+1/2,i}] \\ &\quad + \|\mathbb{E}[V_{t,k+1,i} - s | \mathcal{F}_{t,k+1/2,i}]\|^2. \end{aligned} \quad (39)$$

On one hand,

$$\|V_{t,k+1,i} - \mathbb{E}[V_{t,k+1,i} | \mathcal{F}_{t,k+1/2,i}]\|^2 = \alpha^2 \|\text{Quant}(\Delta_{t,k+1,i}) - \mathbb{E}[\text{Quant}(\Delta_{t,k+1,i}) | \mathcal{F}_{t,k+1/2,i}]\|^2$$

and by A6,

$$\mathbb{E} [\|V_{t,k+1,i} - \mathbb{E}[V_{t,k+1,i} | \mathcal{F}_{t,k+1/2,i}]\|^2 | \mathcal{F}_{t,k+1/2,i}] \leq \alpha^2 \omega \|\Delta_{t,k+1,i}\|^2. \quad (40)$$

On the other hand, for any $s \in \mathbb{R}^q$, and using Lemma 6

$$\begin{aligned} \|\mathbb{E}[V_{t,k+1,i} - s | \mathcal{F}_{t,k+1/2,i}]\|^2 &= \|(1 - \alpha)(V_{t,k,i} - s) + \alpha(S_{t,k+1,i} - \widehat{S}_{t,k} - s)\|^2 \\ &= (1 - \alpha) \|V_{t,k,i} - s\|^2 + \alpha \|S_{t,k+1,i} - \widehat{S}_{t,k} - s\|^2 - \alpha(1 - \alpha) \|V_{t,k,i} - S_{t,k+1,i} + \widehat{S}_{t,k}\|^2 \\ &= (1 - \alpha) \|V_{t,k,i} - s\|^2 + \alpha \|S_{t,k+1,i} - \widehat{S}_{t,k} - s\|^2 - \alpha(1 - \alpha) \|\Delta_{t,k+1,i}\|^2. \end{aligned} \quad (41)$$

Let us combine (38) to (41), the last one being applied with $s \leftarrow \mathbf{h}_i(\widehat{S}_{t,k}) \in \mathcal{F}_{t,k,i}^+ \subseteq \mathcal{F}_{t,k+1/2,i}$. Since

$$\|S_{t,k+1,i} - \widehat{S}_{t,k} - \mathbf{h}_i(\widehat{S}_{t,k})\|^2 = \|S_{t,k+1,i} - \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathsf{T}(\widehat{S}_{t,k})\|^2,$$

we write

$$\begin{aligned} \mathbb{E} [\|V_{t,k+1,i} - \mathbf{h}_i(\widehat{S}_{t,k+1})\|^2 | \mathcal{F}_{t,k}] &\leq (1 + \beta^2) L_i^2 \gamma_{t,k+1}^2 \mathbb{E} [\|H_{t,k+1}\|^2 | \mathcal{F}_{t,k}] \\ &\quad + (1 + \beta^{-2}) \left\{ \alpha^2 \omega \mathbb{E} [\|\Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,k}] + (1 - \alpha) \|V_{t,k,i} - \mathbf{h}_i(\widehat{S}_{t,k})\|^2 \right. \\ &\quad \left. + \alpha \mathbb{E} \left[\|S_{t,k+1,i} - \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathsf{T}(\widehat{S}_{t,k})\|^2 | \mathcal{F}_{t,k} \right] - \alpha(1 - \alpha) \mathbb{E} [\|\Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,k}] \right\}. \end{aligned}$$

Choose $\beta^2 > 0$ such that

$$\beta^{-2} \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } \alpha \geq 2/3 \\ \frac{\alpha}{2(1-\alpha)} & \text{if } \alpha \leq 2/3 \end{cases}$$

This implies that

$$(1 + \beta^{-2})(1 - \alpha) \leq 1 - \frac{\alpha}{2}, \quad 1 + \beta^2 \leq \frac{2}{\alpha}, \quad 1 + \beta^{-2} \leq 2.$$

Hence,

$$\begin{aligned} \mathbb{E} [\|V_{t,k+1,i} - \mathbf{h}_i(\widehat{S}_{t,k+1})\|^2 | \mathcal{F}_{t,k}] &\leq (1 - \alpha/2) \|V_{t,k,i} - \mathbf{h}_i(\widehat{S}_{t,k})\|^2 \\ &\quad + \frac{2}{\alpha} L_i^2 \gamma_{t,k+1}^2 \mathbb{E} [\|H_{t,k+1}\|^2 | \mathcal{F}_{t,k}] + 2\alpha \mathbb{E} \left[\|S_{t,k+1,i} - \frac{1}{m} \sum_{j=1}^m \bar{s}_{ij} \circ \mathsf{T}(\widehat{S}_{t,k})\|^2 | \mathcal{F}_{t,k} \right] \\ &\quad + \alpha(\alpha\omega - 1 + \alpha) \mathbb{E} [\|\Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,k}]; \end{aligned}$$

(in the last equality, we use $1 + \beta^{-2} \geq 1$ since $\alpha\omega - 1 + \alpha \leq 0$). Finally, by using Corollary 18, we have

$$\begin{aligned} \mathbb{E} [\|V_{t,k+1,i} - \mathbf{h}_i(\widehat{S}_{t,k+1})\|^2 | \mathcal{F}_{t,0}] &\leq (1 - \alpha/2) \mathbb{E} [\|V_{t,k,i} - \mathbf{h}_i(\widehat{S}_{t,k})\|^2 | \mathcal{F}_{t,0}] \\ &\quad + \frac{2}{\alpha} L_i^2 \gamma_{t,k+1}^2 \mathbb{E} [\|H_{t,k+1}\|^2 | \mathcal{F}_{t,0}] + 2\alpha \frac{L_i^2}{\mathbf{b}} \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} [\|H_{t,\ell}\|^2 | \mathcal{F}_{t,0}] \\ &\quad + \alpha(\alpha\omega - 1 + \alpha) \mathbb{E} [\|\Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,0}]. \end{aligned}$$

The proof is concluded. \square

9.2.5 Results on the random field $H_{t,k+1}$

Proposition 23 shows that the random field $H_{t,k+1}$ is a biased approximation of the field $h(\widehat{S}_{t,k})$, and this bias is canceled at the beginning of each outer loop. Observe also that the bias exists even when there is no compression: when $\omega = 0$ (so that $\text{Quant}(u) = u$) we have

$$\mathbb{E} [H_{t,k+1} | \mathcal{F}_{t,k}] - h(\widehat{S}_{t,k}) = H_{t,k} - h(\widehat{S}_{t,k-1}) ,$$

and the bias is again canceled at the beginning of each outer loop. Proposition 24 provides an upper bound for the variance and the mean squared error of the random field $H_{t,k+1}$. In the case of no compression ($\omega = 0$) and of a single worker ($n = 1$) so that VR-FedEM is SPIDER-EM, Proposition 24 retrieves the variance and the mean squared error of the random field $H_{t,k+1}$ in SPIDER-EM (see [10, Proposition 13]).

Proposition 23. *Assume A6. For any $t \in [k_{\text{out}}]^*$, $\mathbb{E} [H_{t,2} | \mathcal{F}_{t,0}] - h(\widehat{S}_{t,1}) = \mathbb{E} [H_{t,1} | \mathcal{F}_{t,0}] - h(\widehat{S}_{t,0}) = 0$ and for any $k \in [k_{\text{in}} - 1]^*$,*

$$\begin{aligned} \mathbb{E} [H_{t,k+1} | \mathcal{F}_{t,k}] - h(\widehat{S}_{t,k}) &= H_{t,k} - h(\widehat{S}_{t,k-1}) - n^{-1} \sum_{i=1}^n (\text{Quant}(\Delta_{t,k,i}) - \Delta_{t,k,i}) \\ &= n^{-1} \sum_{i=1}^n \left(\mathbb{E} [S_{t,k+1,i} | \mathcal{F}_{t,k}] - m^{-1} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) \right) . \end{aligned}$$

Proof. Let $t \in [k_{\text{out}}]^*$.

• By definition of $H_{t,1}$ and $\Delta_{t,1,i}$, by A6 and by Lemma 21, we have

$$\begin{aligned} \mathbb{E} [H_{t,1} | \mathcal{F}_{t,0}] &= V_{t,0} + n^{-1} \sum_{i=1}^n \mathbb{E} [\text{Quant}(\Delta_{t,1,i}) | \mathcal{F}_{t,0}] = V_{t,0} + n^{-1} \sum_{i=1}^n \mathbb{E} [\Delta_{t,1,i} | \mathcal{F}_{t,0}] \\ &= V_{t,0} + n^{-1} \sum_{i=1}^n \left(\mathbb{E} [S_{t,1,i} | \mathcal{F}_{t,0}] - \widehat{S}_{t,0} - V_{t,0,i} \right) \\ &= n^{-1} \sum_{i=1}^n \mathbb{E} [S_{t,1,i} | \mathcal{F}_{t,0}] - \widehat{S}_{t,0} . \end{aligned}$$

By Proposition 17 $n^{-1} \sum_{i=1}^n \mathbb{E} [S_{t,1,i} | \mathcal{F}_{t,0}] - \widehat{S}_{t,0} = h(\widehat{S}_{t,0})$.

• Consider the case $k = 1$. We have by definition of $H_{t,2}$

$$\mathbb{E} [H_{t,2} | \mathcal{F}_{t,1}] - h(\widehat{S}_{t,1}) = \frac{1}{n} \sum_{i=1}^n \left(\mathbb{E} [S_{t,2,i} | \mathcal{F}_{t,1}] - m^{-1} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,1}) \right) ;$$

Proposition 17 concludes the proof.

• Let $k \geq 2$. As in the case $k = 0$, we have

$$\begin{aligned} \mathbb{E} [H_{t,k+1} | \mathcal{F}_{t,k}] &= V_{t,k} + n^{-1} \sum_{i=1}^n \mathbb{E} [\text{Quant}(\Delta_{t,k+1,i}) | \mathcal{F}_{t,k}] = V_{t,k} + n^{-1} \sum_{i=1}^n \mathbb{E} [\Delta_{t,k+1,i} | \mathcal{F}_{t,k}] \\ &= V_{t,k} + n^{-1} \sum_{i=1}^n \left(\mathbb{E} [S_{t,k+1,i} | \mathcal{F}_{t,k}] - \widehat{S}_{t,k} - V_{t,k,i} \right) \\ &= n^{-1} \sum_{i=1}^n \mathbb{E} [S_{t,k+1,i} | \mathcal{F}_{t,k}] - \widehat{S}_{t,k} , \end{aligned}$$

so that

$$\mathbb{E} [H_{t,k+1} | \mathcal{F}_{t,k}] - h(\widehat{S}_{t,k}) = n^{-1} \sum_{i=1}^n \left(\mathbb{E} [S_{t,k+1,i} | \mathcal{F}_{t,k}] - m^{-1} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) \right) . \quad (42)$$

By Proposition 17, upon noting that $\mathcal{F}_{t,k} \subset \mathcal{F}_{t,k,i}^+$ and $S_{t,k,i}, \widehat{S}_{t,k-1} \in \mathcal{F}_{t,k}$, we have

$$n^{-1} \sum_{i=1}^n \mathbb{E} [S_{t,k+1,i} | \mathcal{F}_{t,k}] - m^{-1} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k}) = n^{-1} \sum_{i=1}^n \left(S_{t,k,i} - m^{-1} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k-1}) \right). \quad (43)$$

On the other hand, observe that

$$\begin{aligned} H_{t,k} &= V_{t,k-1} + n^{-1} \sum_{i=1}^n \text{Quant}(\Delta_{t,k,i}) \\ &= V_{t,k-1} + n^{-1} \sum_{i=1}^n S_{t,k,i} - \widehat{S}_{t,k-1} - n^{-1} \sum_{i=1}^n V_{t,k-1,i} + n^{-1} \sum_{i=1}^n (\text{Quant}(\Delta_{t,k,i}) - \Delta_{t,k,i}) \\ &= n^{-1} \sum_{i=1}^n S_{t,k,i} - \widehat{S}_{t,k-1} + n^{-1} \sum_{i=1}^n (\text{Quant}(\Delta_{t,k,i}) - \Delta_{t,k,i}), \end{aligned}$$

where we used Lemma 21. This yields

$$\begin{aligned} H_{t,k} - h(\widehat{S}_{t,k-1}) &= n^{-1} \sum_{i=1}^n \left(S_{t,k,i} - m^{-1} \sum_{j=1}^m \bar{s}_{ij} \circ \mathbb{T}(\widehat{S}_{t,k-1}) \right) + n^{-1} \sum_{i=1}^n (\text{Quant}(\Delta_{t,k,i}) - \Delta_{t,k,i}). \quad (44) \end{aligned}$$

The proof is concluded by combining (42), (43) and (44). \square

Proposition 24. *Assume A6 and A9. For any $t \in [k_{\text{out}}]^*$,*

$$\mathbb{E} \left[\|H_{t,1} - h(\widehat{S}_{t,0})\|^2 | \mathcal{F}_{t,0} \right] \leq \frac{\omega}{n} \left(\frac{1}{n} \sum_{i=1}^n \|V_{t,0,i} - h_i(\widehat{S}_{t,0})\|^2 \right),$$

and for any $k \in [k_{\text{in}} - 1]^*$,

$$\begin{aligned} \mathbb{E} \left[\|H_{t,k+1} - h(\widehat{S}_{t,k})\|^2 | \mathcal{F}_{t,0} \right] &\leq \frac{\omega}{n} \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\|\Delta_{t,k+1,i}\|^2 | \mathcal{F}_{t,0} \right] + \frac{L^2}{nb} \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} \left[\|H_{t,\ell}\|^2 | \mathcal{F}_{t,0} \right], \\ \mathbb{E} \left[\|\mathbb{E} [H_{t,k+1} | \mathcal{F}_{t,k}] - h(\widehat{S}_{t,k})\|^2 | \mathcal{F}_{t,0} \right] &\leq \frac{L^2}{nb} \sum_{\ell=1}^{k-1} \gamma_{t,\ell}^2 \mathbb{E} \left[\|H_{t,\ell}\|^2 | \mathcal{F}_{t,0} \right]. \end{aligned}$$

Proof. • **Case $k = 1$.** From Proposition 23 and the definition of $H_{t,1}$, we have

$$\begin{aligned} H_{t,1} - h(\widehat{S}_{t,0}) &= H_{t,1} - \mathbb{E} [H_{t,1} | \mathcal{F}_{t,0}] = n^{-1} \sum_{i=1}^n (\text{Quant}(\Delta_{t,1,i}) - \mathbb{E} [\text{Quant}(\Delta_{t,1,i}) | \mathcal{F}_{t,0}]) \\ &= n^{-1} \sum_{i=1}^n (\text{Quant}(\Delta_{t,1,i}) - \mathbb{E} [\Delta_{t,1,i} | \mathcal{F}_{t,0}]), \end{aligned}$$

where we used $\mathbb{E} [\text{Quant}(\Delta_{t,1,i}) | \mathcal{F}_{t,1/2,i}] = \Delta_{t,1,i}$ and $\mathcal{F}_{t,0} \subset \mathcal{F}_{t,1/2,i}$ in the last equality. In addition, since $\widehat{S}_{t,0} = \widehat{S}_{t,-1}$, we have (see Proposition 17)

$$S_{t,1,i} = S_{t,0,i} = h_i(\widehat{S}_{t,0}) + \widehat{S}_{t,0}.$$

Hence,

$$\Delta_{t,1,i} = S_{t,1,i} - \widehat{S}_{t,0} - V_{t,0,i} = h_i(\widehat{S}_{t,0}) - V_{t,0,i}.$$

Therefore, $\mathbb{E} [\Delta_{t,1,i} | \mathcal{F}_{t,0}] = \Delta_{t,1,i} = h_i(\widehat{S}_{t,0}) - V_{t,0,i}$. Since the workers are independent, we write

$$\mathbb{E} \left[\|H_{t,1} - h(\widehat{S}_{t,0})\|^2 | \mathcal{F}_{t,0} \right] = \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[\|\text{Quant}(h_i(\widehat{S}_{t,0}) - V_{t,0,i}) - (h_i(\widehat{S}_{t,0}) - V_{t,0,i})\|^2 | \mathcal{F}_{t,0} \right].$$

By A6, this yields

$$\mathbb{E} \left[\|H_{t,1} - \mathfrak{h}(\widehat{S}_{t,0})\|^2 \middle| \mathcal{F}_{t,0} \right] \leq \frac{\omega}{n} \frac{1}{n} \sum_{i=1}^n \|\mathfrak{h}_i(\widehat{S}_{t,0}) - V_{t,0,i}\|^2.$$

• **Case $k \geq 1$.** Let $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}} - 1]^*$. We write

$$\begin{aligned} \mathbb{E} \left[\|H_{t,k+1} - \mathfrak{h}(\widehat{S}_{t,k})\|^2 \middle| \mathcal{F}_{t,0} \right] &= \mathbb{E} \left[\|H_{t,k+1} - \mathbb{E}[H_{t,k+1} | \mathcal{F}_{t,k}]\|^2 \middle| \mathcal{F}_{t,0} \right] \\ &\quad + \mathbb{E} \left[\|\mathbb{E}[H_{t,k+1} | \mathcal{F}_{t,k}] - \mathfrak{h}(\widehat{S}_{t,k})\|^2 \middle| \mathcal{F}_{t,0} \right]. \end{aligned} \quad (45)$$

Let us first consider the bias term. From Proposition 17, Proposition 23 and the definition of $S_{t,k+1,i}$ (remember that $S_{t,k,i}$, $\widehat{S}_{t,k}$ and $\widehat{S}_{t,k-1}$ are in $\mathcal{F}_{t,k,i}^+ \supset \mathcal{F}_{t,k}$), it holds

$$\begin{aligned} &\mathbb{E} \left[\left\| \mathbb{E}[H_{t,k+1} | \mathcal{F}_{t,k}] - \mathfrak{h}(\widehat{S}_{t,k}) \right\|^2 \middle| \mathcal{F}_{t,0} \right] \\ &= \mathbb{E} \left[\left\| n^{-1} \sum_{i=1}^n (\mathbb{E}[S_{t,k+1,i} | \mathcal{F}_{t,k}] - m^{-1} \sum_{j=1}^m \bar{s}_{ij} \circ \mathfrak{T}(\widehat{S}_{t,k})) \right\|^2 \middle| \mathcal{F}_{t,0} \right] \\ &\leq \mathbb{E} \left[\left\| n^{-1} \sum_{i=1}^n (S_{t,k,i} - m^{-1} \sum_{j=1}^m \bar{s}_{ij} \circ \mathfrak{T}(\widehat{S}_{t,k-1})) \right\|^2 \middle| \mathcal{F}_{t,0} \right]. \end{aligned}$$

By Proposition 17 again, the RHS is zero when $k = 1$; when $k \geq 2$, by Corollary 18 and the independence of the workers, we have yields

$$\mathbb{E} \left[\|\mathbb{E}[H_{t,k+1} | \mathcal{F}_{t,k}] - \mathfrak{h}(\widehat{S}_{t,k})\|^2 \middle| \mathcal{F}_{t,0} \right] \leq \frac{L^2}{nb} \sum_{\ell=1}^{k-1} \gamma_{t,\ell}^2 \mathbb{E} [\|H_{t,\ell}\|^2 | \mathcal{F}_{t,0}]. \quad (46)$$

Let us now consider the variance term. We have from the definition of $H_{t,k+1}$ and A6

$$H_{t,k+1} - \mathbb{E}[H_{t,k+1} | \mathcal{F}_{t,k}] = \frac{1}{n} \sum_{i=1}^n (\text{Quant}(\Delta_{t,k+1,i}) - \mathbb{E}[\Delta_{t,k+1,i} | \mathcal{F}_{t,k}])$$

and here again, by the independence of the workers

$$\begin{aligned} &\mathbb{E} \left[\|H_{t,k+1} - \mathbb{E}[H_{t,k+1} | \mathcal{F}_{t,k}]\|^2 \middle| \mathcal{F}_{t,0} \right] \\ &\leq \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[\|\text{Quant}(\Delta_{t,k+1,i}) - \mathbb{E}[\Delta_{t,k+1,i} | \mathcal{F}_{t,k}]\|^2 \middle| \mathcal{F}_{t,0} \right]. \end{aligned} \quad (47)$$

The proof follows from (45) to (47) and Proposition 20. \square

9.3 Proof of Theorem 4

Theorem 4 is a corollary of the more general following proposition.

Proposition 25. *Assume A1 to 3, A4, A6 and A9. Set $L^2 \stackrel{\text{def}}{=} n^{-1}m^{-1} \sum_{i=1}^n \sum_{j=1}^m L_{ij}^2$. Let $\{\widehat{S}_{t,k}, t \in [k_{\text{out}}]^*, k \in [k_{\text{in}} - 1]\}$ be given by algorithm 2 run with any $\alpha \leq 1/(1 + \omega)$, and $b \geq 1$, with $V_{1,0,i} = \mathfrak{h}_i(\widehat{S}_{1,0})$ for any $i \in [n]^*$. Let (τ, K) be a uniform random variable on $[k_{\text{out}}]^* \times [k_{\text{in}} - 1]$, independent of $\{\widehat{S}_{t,k}, t \in [k_{\text{out}}]^*, k \in [k_{\text{in}} - 1]\}$. Then, it holds*

$$v_{\min} (1 - \gamma \Lambda_\star) \mathbb{E} [\|H_{\tau,K+1}\|^2] \leq \gamma^{-1} k_{\text{in}}^{-1} k_{\text{out}}^{-1} \left(\mathbb{E} [\mathfrak{W}(\widehat{S}_{1,\star})] - \min \mathfrak{W} \right),$$

where

$$\Lambda_\star \stackrel{\text{def}}{=} \frac{L_{\dot{\mathfrak{W}}}}{2v_{\min}} + 2\sqrt{2} \frac{v_{\max}}{v_{\min}} \frac{L}{\sqrt{n}\alpha} \left(\omega + \frac{k_{\text{in}}\alpha^2}{8b} (1 + 10\omega) \right)^{1/2}.$$

The proof of Theorem 4 from Proposition 25 (which corresponds to particular choices of b, α , etc. is detailed in Section 9.4).

9.3.1 Control of $H_{\tau,K}$

Let $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}} - 1]$. By A4, we have

$$\mathbb{W}(\widehat{S}_{t,k+1}) \leq \mathbb{W}(\widehat{S}_{t,k}) + \left\langle \nabla \mathbb{W}(\widehat{S}_{t,k}), \widehat{S}_{t,k+1} - \widehat{S}_{t,k} \right\rangle + \frac{L_{\dot{\mathbb{W}}}}{2} \|\widehat{S}_{t,k+1} - \widehat{S}_{t,k}\|^2.$$

Since $\widehat{S}_{t,k+1} - \widehat{S}_{t,k} = \gamma_{t,k+1} H_{t,k+1}$, we have using again A4

$$\mathbb{W}(\widehat{S}_{t,k+1}) \leq \mathbb{W}(\widehat{S}_{t,k}) - \gamma_{t,k+1} \left\langle B(\widehat{S}_{t,k}) \mathbf{h}(\widehat{S}_{t,k}), H_{t,k+1} \right\rangle + \frac{L_{\dot{\mathbb{W}}}}{2} \gamma_{t,k+1}^2 \|H_{t,k+1}\|^2.$$

We have the inequality, for any $\beta > 0$:

$$-\langle Bh, H \rangle \leq -\langle BH, H \rangle - \langle B(h - H), H \rangle \leq -\langle BH, H \rangle + \frac{\beta^2}{2} \|H\|^2 + \frac{1}{2\beta^2} \|B(H - h)\|^2.$$

By A4 again, this inequality yields for any $\beta_{t,k+1} > 0$ after applying the conditional expectation

$$\begin{aligned} \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k+1}) \middle| \mathcal{F}_{t,0} \right] &\leq \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k}) \middle| \mathcal{F}_{t,0} \right] - \gamma_{t,k+1} v_{\min} \Lambda_{t,k+1} \mathbb{E} \left[\|H_{t,k+1}\|^2 \middle| \mathcal{F}_{t,0} \right] \\ &\quad + \frac{\gamma_{t,k+1}}{2\beta_{t,k+1}^2} v_{\max}^2 \mathbb{E} \left[\|H_{t,k+1} - \mathbf{h}(\widehat{S}_{t,k})\|^2 \middle| \mathcal{F}_{t,0} \right], \end{aligned} \quad (48)$$

where

$$\Lambda_{t,k+1} \stackrel{\text{def}}{=} 1 - \gamma_{t,k+1} \frac{L_{\dot{\mathbb{W}}}}{2v_{\min}} - \frac{\beta_{t,k+1}^2}{2v_{\min}}.$$

By (48) and Proposition 24, it holds

$$\begin{aligned} \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k+1}) \middle| \mathcal{F}_{t,0} \right] &\leq \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k}) \middle| \mathcal{F}_{t,0} \right] - \gamma_{t,k+1} v_{\min} \Lambda_{t,k+1} \mathbb{E} \left[\|H_{t,k+1}\|^2 \middle| \mathcal{F}_{t,0} \right] \\ &\quad + \frac{\gamma_{t,k+1}}{2\beta_{t,k+1}^2} v_{\max}^2 \frac{L^2}{n\mathbf{b}} \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} \left[\|H_{t,\ell}\|^2 \middle| \mathcal{F}_{t,0} \right] \\ &\quad + \frac{\gamma_{t,k+1}}{2\beta_{t,k+1}^2} v_{\max}^2 \frac{\omega}{n} \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\|\Delta_{t,k+1,i}\|^2 \middle| \mathcal{F}_{t,0} \right]. \end{aligned} \quad (49)$$

Set

$$G_{t,k} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|V_{t,k,i} - \mathbf{h}_i(\widehat{S}_{t,k})\|^2.$$

From Proposition 19, we obtain

$$\begin{aligned} \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k+1}) \middle| \mathcal{F}_{t,0} \right] &\leq \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k}) \middle| \mathcal{F}_{t,0} \right] - \gamma_{t,k+1} v_{\min} \Lambda_{t,k+1} \mathbb{E} \left[\|H_{t,k+1}\|^2 \middle| \mathcal{F}_{t,0} \right] \\ &\quad + \frac{\gamma_{t,k+1}}{2\beta_{t,k+1}^2} v_{\max}^2 \frac{L^2}{n\mathbf{b}} (1 + 2\omega) \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} \left[\|H_{t,\ell}\|^2 \middle| \mathcal{F}_{t,0} \right] + \frac{\gamma_{t,k+1}}{\beta_{t,k+1}^2} v_{\max}^2 \frac{\omega}{n} \mathbb{E} \left[G_{t,k} \middle| \mathcal{F}_{t,0} \right]. \end{aligned} \quad (50)$$

Assume that $k \mapsto \gamma_{t,k+1}/\beta_{t,k+1}^2$ is a non-increasing sequence and set

$$C_{t,k+1} \stackrel{\text{def}}{=} \frac{2\omega}{\alpha n} v_{\max}^2 \frac{\gamma_{t,k+1}}{\beta_{t,k+1}^2}. \quad (51)$$

From Proposition 22, since $\alpha \in (0, 1/(1 + \omega)]$, we have

$$\begin{aligned} C_{t,k+1} \mathbb{E} \left[G_{t,k+1} \middle| \mathcal{F}_{t,0} \right] &\leq (1 - \alpha/2) C_{t,k+1} \mathbb{E} \left[G_{t,k} \middle| \mathcal{F}_{t,0} \right] + \frac{2}{\alpha} L^2 \gamma_{t,k+1}^2 C_{t,k+1} \mathbb{E} \left[\|H_{t,k+1}\|^2 \middle| \mathcal{F}_{t,0} \right] \\ &\quad + 2\alpha \frac{L^2}{\mathbf{b}} C_{t,k+1} \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} \left[\|H_{t,\ell}\|^2 \middle| \mathcal{F}_{t,0} \right]. \end{aligned} \quad (52)$$

Upon noting that by definition of $C_{t,k+1}$ we have (remember that $C_{t,k+1} \leq C_{t,k}$)

$$(1 - \alpha/2)C_{t,k+1} - C_{t,k} + \frac{\gamma_{t,k+1}}{\beta_{t,k+1}^2} v_{\max}^2 \frac{\omega}{n} \leq 0 ,$$

this yields from (50) and (52)

$$\begin{aligned} \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k+1}) | \mathcal{F}_{t,0} \right] + C_{t,k+1} \mathbb{E} [G_{t,k+1} | \mathcal{F}_{t,0}] &\leq \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k}) | \mathcal{F}_{t,0} \right] + C_{t,k} \mathbb{E} [G_{t,k} | \mathcal{F}_{t,0}] \\ &- \left(\gamma_{t,k+1} v_{\min} \Lambda_{t,k+1} - \frac{2}{\alpha} L^2 \gamma_{t,k+1}^2 C_{t,k+1} \right) \mathbb{E} [\|H_{t,k+1}\|^2 | \mathcal{F}_{t,0}] \\ &+ \left(\frac{\gamma_{t,k+1}}{2\beta_{t,k+1}^2} v_{\max}^2 \frac{L^2}{nb} (1 + 2\omega) + 2\alpha \frac{L^2}{b} C_{t,k+1} \right) \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} [\|H_{t,\ell}\|^2 | \mathcal{F}_{t,0}] . \end{aligned}$$

Let us restrict the computations to the case $\gamma_{t,k} = \gamma$, $\beta_{t,k} = \beta$ (which implies $C_{t,k+1} = C_{t,k} =: C$); we obtain

$$\begin{aligned} \gamma v_{\min} \left(1 - \gamma \frac{L\dot{\mathbb{W}}}{2v_{\min}} - \frac{\beta^2}{2v_{\min}} - \frac{\gamma^2}{\beta^2} \frac{4v_{\max}^2}{v_{\min}} L^2 \frac{\omega}{\alpha^2 n} \right) \mathbb{E} [\|H_{t,k+1}\|^2 | \mathcal{F}_{t,0}] \\ \leq \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k}) | \mathcal{F}_{t,0} \right] + C \mathbb{E} [G_{t,k} | \mathcal{F}_{t,0}] - \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k+1}) | \mathcal{F}_{t,0} \right] - C \mathbb{E} [G_{t,k+1} | \mathcal{F}_{t,0}] \\ + \frac{\gamma^3}{2\beta^2} v_{\max}^2 \frac{L^2}{nb} (1 + 10\omega) \sum_{\ell=1}^k \mathbb{E} [\|H_{t,\ell}\|^2 | \mathcal{F}_{t,0}] . \end{aligned}$$

We now sum from $k = 0$ to $k = k_{\text{in}} - 1$ and divide by k_{in} :

$$\begin{aligned} \gamma v_{\min} \left(1 - \gamma \frac{L\dot{\mathbb{W}}}{2v_{\min}} - \frac{\beta^2}{2v_{\min}} - \frac{\gamma^2}{\beta^2} \frac{4v_{\max}^2}{v_{\min}} L^2 \frac{\omega}{\alpha^2 n} \right) \frac{1}{k_{\text{in}}} \sum_{k=1}^{k_{\text{in}}} \mathbb{E} [\|H_{t,k}\|^2 | \mathcal{F}_{t,0}] \\ \leq k_{\text{in}}^{-1} \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,0}) | \mathcal{F}_{t,0} \right] + \frac{C}{k_{\text{in}}} \mathbb{E} [G_{t,0} | \mathcal{F}_{t,0}] \\ - k_{\text{in}}^{-1} \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k_{\text{in}}}) | \mathcal{F}_{t,0} \right] - \frac{C}{k_{\text{in}}} \mathbb{E} [G_{t,k_{\text{in}}} | \mathcal{F}_{t,0}] \\ + \frac{\gamma^3}{2\beta^2} v_{\max}^2 \frac{L^2}{nb} (1 + 10\omega) \sum_{k=1}^{k_{\text{in}}} \mathbb{E} [\|H_{t,k}\|^2 | \mathcal{F}_{t,0}] . \end{aligned}$$

As a conclusion, we have

$$\begin{aligned} \gamma v_{\min} \left(1 - \gamma \frac{L\dot{\mathbb{W}}}{2v_{\min}} - \gamma \bar{\Lambda} \right) \frac{1}{k_{\text{in}}} \sum_{k=0}^{k_{\text{in}}-1} \mathbb{E} [\|H_{t,k+1}\|^2 | \mathcal{F}_{t,0}] \\ \leq k_{\text{in}}^{-1} \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,0}) | \mathcal{F}_{t,0} \right] + \frac{C}{k_{\text{in}}} \mathbb{E} [G_{t,0} | \mathcal{F}_{t,0}] \\ - k_{\text{in}}^{-1} \mathbb{E} \left[\mathbb{W}(\widehat{S}_{t,k_{\text{in}}}) | \mathcal{F}_{t,0} \right] - \frac{C}{k_{\text{in}}} \mathbb{E} [G_{t,k_{\text{in}}} | \mathcal{F}_{t,0}] . \end{aligned}$$

where

$$\bar{\Lambda} \stackrel{\text{def}}{=} \frac{\beta^2}{2v_{\min}\gamma} + \frac{\gamma}{\beta^2} \frac{4v_{\max}^2}{v_{\min}} L^2 \frac{\omega}{\alpha^2 n} + \frac{\gamma}{2\beta^2} \frac{v_{\max}^2}{v_{\min}} \frac{L^2 k_{\text{in}}}{nb} (1 + 10\omega) .$$

Next, we sum from $t = 1$ to $t = k_{\text{out}}$, divide by k_{out} .

$$\begin{aligned} \gamma v_{\min} \left(1 - \gamma \frac{L\dot{\mathbb{W}}}{2v_{\min}} - \gamma \bar{\Lambda} \right) \frac{1}{k_{\text{out}} k_{\text{in}}} \sum_{k=1}^{k_{\text{out}}} \sum_{k=1}^{k_{\text{in}}} \mathbb{E} [\|H_{t,k+1}\|^2] \\ \leq k_{\text{in}}^{-1} k_{\text{out}}^{-1} \left(\mathbb{E} \left[\mathbb{W}(\widehat{S}_{1,0}) \right] - \min \mathbb{W} \right) + \frac{C}{k_{\text{in}} k_{\text{out}}} \mathbb{E} [G_{1,0}] . \quad (53) \end{aligned}$$

Finally, we apply the expectation, with (τ, K) a uniform random variable on $[k_{\text{out}}]^* \times [k_{\text{in}} - 1]$, independent of $\{\widehat{S}_{t,k}, t \in [k_{\text{out}}]^*, k \in [k_{\text{in}} - 1]\}$, upon noting that $G_{t,k_{\text{in}}} = G_{t+1,0}$ and $\widehat{S}_{t,k_{\text{in}}} = \widehat{S}_{t+1,0}$, this yields

$$\begin{aligned} \gamma v_{\min} \left(1 - \gamma \frac{L\dot{W}}{2v_{\min}} - \gamma \bar{\Lambda}\right) \mathbb{E} [\|H_{\tau, K+1}\|^2] \\ \leq k_{\text{in}}^{-1} k_{\text{out}}^{-1} \left(\mathbb{E} [\mathbb{W}(\widehat{S}_{1,0})] - \min \mathbb{W}\right) + \frac{C}{k_{\text{in}} k_{\text{out}}} \mathbb{E} [G_{1,0}]. \end{aligned} \quad (54)$$

Impact of initialization. With $V_{1,0,i} = h_i(\widehat{S}_{1,0})$ for any $i \in [n]^*$, we have $G_{1,0} = 0$.

Choice of β . The latter inequality is true for all parameter $\beta^2 > 0$ (coming from Young's inequality). We can thus optimize the value of β^2 to minimize the value of $\bar{\Lambda}$. We here discuss this choice. First, to ensure that $\bar{\Lambda}$ is independent of γ , we introduce \mathbf{a} , and set $\beta^2 = \mathbf{a}\gamma$ so that

$$\begin{aligned} \bar{\Lambda} &= \frac{\mathbf{a}}{2v_{\min}} + \frac{1}{\mathbf{a}} \frac{4v_{\max}^2 L^2}{v_{\min}} \frac{\omega}{\alpha^2 n} + \frac{1}{2\mathbf{a}} \frac{v_{\max}^2 L^2 k_{\text{in}}}{v_{\min} n \mathbf{b}} (1 + 10\omega) \\ &= \frac{\mathbf{a}}{2v_{\min}} + \frac{4}{\mathbf{a}} \frac{v_{\max}^2 L^2}{v_{\min} n \alpha^2} \left(\omega + \frac{k_{\text{in}} \alpha^2}{8\mathbf{b}} (1 + 10\omega)\right). \end{aligned}$$

Next, we optimize the value of \mathbf{a} .² Upon noting that $\mathbf{a} \mapsto A\mathbf{a} + B/\mathbf{a}$ (for $A, B > 0$) is lower bounded by $2\sqrt{AB}$ and its minimizer is $\mathbf{a}_* \stackrel{\text{def}}{=} \sqrt{B/A}$, we choose

$$\mathbf{a}_* \stackrel{\text{def}}{=} 2\sqrt{2} v_{\max} \frac{L}{\sqrt{n}\alpha} \left(\omega + \frac{k_{\text{in}} \alpha^2}{8\mathbf{b}} (1 + 10\omega)\right)^{1/2}$$

and obtain

$$\bar{\Lambda} = 2\sqrt{2} \frac{v_{\max}}{v_{\min}} \frac{L}{\sqrt{n}\alpha} \left(\omega + \frac{k_{\text{in}} \alpha^2}{8\mathbf{b}} (1 + 10\omega)\right)^{1/2}. \quad (55)$$

Combining Equation (55) and Equation (54), we obtain

$$v_{\min} (1 - \gamma \Lambda_*) \mathbb{E} [\|H_{\tau, K+1}\|^2] \leq \gamma^{-1} k_{\text{in}}^{-1} k_{\text{out}}^{-1} \left(\mathbb{E} [\mathbb{W}(\widehat{S}_{1,0})] - \min \mathbb{W}\right),$$

where

$$\Lambda_* \stackrel{\text{def}}{=} \frac{L\dot{W}}{2v_{\min}} + 2\sqrt{2} \frac{v_{\max}}{v_{\min}} \frac{L}{\sqrt{n}\alpha} \left(\omega + \frac{k_{\text{in}} \alpha^2}{8\mathbf{b}} (1 + 10\omega)\right)^{1/2}, \quad (56)$$

which is the result of Proposition 25.

9.4 Proof of Theorem 4 (Equation (13)) from Proposition 25

We apply Proposition 25 with: $\mathbf{b} \stackrel{\text{def}}{=} \lceil \frac{k_{\text{in}}}{(1+\omega)^2} \rceil$ and the largest possible learning rate $\alpha = (1 + \omega)^{-1}$: this gives in Equation (56)

$$\begin{aligned} \Lambda_* &= \frac{L\dot{W}}{2v_{\min}} + 2\sqrt{2} \frac{v_{\max}}{v_{\min}} \frac{L}{\sqrt{n}} (1 + \omega) \left(\omega + \frac{1 + 10\omega}{8}\right)^{1/2} \\ &= \frac{L\dot{W}}{2v_{\min}} \left(1 + 4\sqrt{2} \frac{v_{\max}}{L\dot{W}} \frac{L}{\sqrt{n}} (1 + \omega) \left(\omega + \frac{1 + 10\omega}{8}\right)^{1/2}\right). \end{aligned}$$

Next, we choose γ to be the largest possible value to ensure $(1 - \gamma \Lambda_*) \geq \frac{1}{2}$. For all t, k ,

$$\gamma_{t,k} = \gamma \stackrel{\text{def}}{=} \frac{1}{2\Lambda_*} = \frac{v_{\min}}{L\dot{W}} \left(1 + 4\sqrt{2} \frac{v_{\max}}{L\dot{W}} \frac{L}{\sqrt{n}} (1 + \omega) \left(\omega + \frac{1 + 10\omega}{8}\right)^{1/2}\right)^{-1}.$$

This gives the first result of Theorem 4, namely Equation (13). We give the proof of the second result, Equation (14) in the following subsection.

²Remark that this optimization step is crucial to optimize the dependency of $\bar{\Lambda}$ w.r.t. ω : this ensures that $\bar{\Lambda} \propto \omega^{3/2}$.

9.5 Proof of Theorem 4 (Equation (14)): control on $h(\widehat{S}_{\tau,K})$

We now establish (14) for $\gamma_{t,k} = \gamma$. Let $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}} - 1]$. We have

$$\|h(\widehat{S}_{t,k})\|^2 \leq 2\mathbb{E} [H_{t,k+1}|\mathcal{F}_{t,k}]^2 + 2\|h(\widehat{S}_{t,k}) - \mathbb{E} [H_{t,k+1}|\mathcal{F}_{t,k}]\|^2. \quad (57)$$

Let us consider the first term in (57). By Jensen's inequality and the tower property of conditional expectations

$$\mathbb{E} [\|\mathbb{E} [H_{t,k+1}|\mathcal{F}_{t,k}]\|^2|\mathcal{F}_{t,0}] \leq \mathbb{E} [\mathbb{E} [\|H_{t,k+1}\|^2|\mathcal{F}_{t,k}]|\mathcal{F}_{t,0}] = \mathbb{E} [\|H_{t,k+1}\|^2|\mathcal{F}_{t,0}].$$

Let us now consider the second term in (57). By Proposition 23 and Proposition 24, we have

$$\mathbb{E} [\|\mathbb{E} [H_{t,k+1}|\mathcal{F}_{t,k}] - h(\widehat{S}_{t,k})\|^2|\mathcal{F}_{t,0}] \leq \begin{cases} \gamma^2 \frac{L^2}{nb} \sum_{\ell=1}^{k-1} \mathbb{E} [\|H_{t,\ell}\|^2|\mathcal{F}_{t,0}] & \text{when } k \geq 2 \\ 0 & \text{when } k \in \{0, 1\} \end{cases}.$$

Therefore, we write

$$\mathbb{E} [\|h(\widehat{S}_{t,k})\|^2] \leq 2\mathbb{E} [\|H_{t,k+1}\|^2] + 2\gamma^2 \frac{L^2}{nb} \sum_{\ell=1}^{k-1} \mathbb{E} [\|H_{t,\ell}\|^2]$$

We now sum from $k = 0$ to $k = k_{\text{in}} - 1$, then from $t = 1$ to $t = k_{\text{out}}$, and finally we divide by $k_{\text{in}}k_{\text{out}}$. This yields

$$\begin{aligned} \mathbb{E} [\|h(\widehat{S}_{\tau,K})\|^2] &\leq 2\mathbb{E} [\|H_{\tau,K+1}\|^2] + 2\gamma^2 \frac{L^2}{nb} \frac{1}{k_{\text{in}}k_{\text{out}}} \sum_{t=1}^{k_{\text{out}}} \sum_{k=2}^{k_{\text{in}}-1} \sum_{\ell=1}^{k-1} \mathbb{E} [\|H_{t,\ell}\|^2] \\ &\leq 2\mathbb{E} [\|H_{\tau,K+1}\|^2] + 2\gamma^2 \frac{L^2}{nb} \frac{1}{k_{\text{out}}} \sum_{t=1}^{k_{\text{out}}} \sum_{k=1}^{k_{\text{in}}-2} \mathbb{E} [\|H_{t,k}\|^2] \\ &\leq 2\mathbb{E} [\|H_{\tau,K+1}\|^2] + 2\gamma^2 \frac{L^2}{n} \frac{k_{\text{in}}}{b} \mathbb{E} [\|H_{\tau,K+1}\|^2] \\ &\leq 2 \left(1 + \gamma^2 \frac{L^2}{n} \frac{k_{\text{in}}}{b} \right) \mathbb{E} [\|H_{\tau,K+1}\|^2]. \end{aligned}$$

9.6 On the convergence of the $V_{t,k,i}$'s

In this subsection, we provide a complementary result, to support the assertion made in the text, that the variable $V_{t,k,i}$ approximates $h_i(\widehat{S}_{t,k})$. Recall that for $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}}]$, $G_{t,k} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \|V_{t,k,i} - h_i(\widehat{S}_{t,k})\|^2$.

Proposition 26. *When running algorithm 2 with a constant step size γ equal to*

$$\gamma \stackrel{\text{def}}{=} \frac{v_{\min}}{L_{\widehat{W}}} \left(1 + 4\sqrt{2} \frac{v_{\max}}{L_{\widehat{W}}} \frac{L}{\sqrt{n}} (1 + \omega) \left(\omega + \frac{1 + 10\omega}{8} \right)^{1/2} \right)^{-1},$$

with $b \stackrel{\text{def}}{=} \lceil \frac{k_{\text{in}}}{(1+\omega)^2} \rceil$ and $\alpha \stackrel{\text{def}}{=} 1/(\omega + 1)$, we have

$$\frac{1}{k_{\text{out}}k_{\text{in}}} \sum_{t=1}^{k_{\text{out}}} \sum_{k=1}^{k_{\text{in}}} \mathbb{E}[G_{t,k}] \leq \frac{2(1+\omega)}{k_{\text{in}}k_{\text{out}}} \mathbb{E}[G_{1,0}] + 16 \frac{\gamma}{k_{\text{in}}k_{\text{out}}} \frac{(1+\omega)^2 L^2}{v_{\min}} \left(\mathbb{E} [\mathbb{W}(\widehat{S}_{1,0})] - \min \mathbb{W} \right).$$

In words, the Cesaro average $\frac{1}{k_{\text{out}}k_{\text{in}}} \sum_{t=1}^{k_{\text{out}}} \sum_{k=1}^{k_{\text{in}}} \mathbb{E}[G_{t,k}]$ decreases proportionally to the number of iterations $k_{\text{in}}k_{\text{out}}$. Consequently, the average squared distance between $V_{t,k,i}$ and $h_i(\widehat{S}_{t,k})$ (i.e., $G_{t,k}$), converges to 0 in the sense of Cesaro.

Proof. From Proposition 22, we have that, $t \in [k_{\text{out}}]^*$ and $k \in [k_{\text{in}}]$, and any $\alpha \leq (\omega + 1)^{-1}$:

$$\begin{aligned} \mathbb{E} [G_{t,k+1}|\mathcal{F}_{t,0}] &\leq (1 - \alpha/2) \mathbb{E} [G_{t,k}|\mathcal{F}_{t,0}] \\ &\quad + \frac{2}{\alpha} L^2 \gamma_{t,k+1}^2 \mathbb{E} [\|H_{t,k+1}\|^2|\mathcal{F}_{t,0}] + 2\alpha \frac{L^2}{b} \sum_{\ell=1}^k \gamma_{t,\ell}^2 \mathbb{E} [\|H_{t,\ell}\|^2|\mathcal{F}_{t,0}]. \end{aligned}$$

Equivalently:

$$\begin{aligned} \alpha/2\mathbb{E}[G_{t,k}|\mathcal{F}_{t,0}] &\leq \mathbb{E}[G_{t,k}|\mathcal{F}_{t,0}] - \mathbb{E}[G_{t,k+1}|\mathcal{F}_{t,0}] \\ &\quad + \frac{2}{\alpha}L^2\gamma_{t,k+1}^2\mathbb{E}[\|H_{t,k+1}\|^2|\mathcal{F}_{t,0}] + 2\alpha\frac{L^2}{\mathbf{b}}\sum_{\ell=1}^k\gamma_{t,\ell}^2\mathbb{E}[\|H_{t,\ell}\|^2|\mathcal{F}_{t,0}]. \end{aligned}$$

Summing from $k = 0$ to $k = k_{\text{in}} - 1$, we get, with $\gamma_{t,k+1}^2 = \gamma$:

$$\begin{aligned} \frac{\alpha}{2}\sum_{k=0}^{k_{\text{in}}-1}\mathbb{E}[G_{t,k}|\mathcal{F}_{t,0}] &\leq \mathbb{E}[G_{t,0}|\mathcal{F}_{t,0}] - \mathbb{E}[G_{t,k_{\text{in}}}| \mathcal{F}_{t,0}] \\ &\quad + \frac{2}{\alpha}L^2\gamma^2\sum_{k=1}^{k_{\text{in}}}\mathbb{E}[\|H_{t,k}\|^2|\mathcal{F}_{t,0}] + 2\alpha\frac{L^2}{\mathbf{b}}\sum_{k=1}^{k_{\text{in}}-1}\sum_{\ell=1}^k\gamma^2\mathbb{E}[\|H_{t,\ell}\|^2|\mathcal{F}_{t,0}] \\ &\leq \mathbb{E}[G_{t,0}|\mathcal{F}_{t,0}] - \mathbb{E}[G_{t,k_{\text{in}}}| \mathcal{F}_{t,0}] \\ &\quad + \frac{2}{\alpha}L^2\gamma^2\sum_{k=1}^{k_{\text{in}}}\mathbb{E}[\|H_{t,k}\|^2|\mathcal{F}_{t,0}] + 2\alpha\frac{L^2k_{\text{in}}}{\mathbf{b}}\sum_{k=1}^{k_{\text{in}}}\gamma^2\mathbb{E}[\|H_{t,k}\|^2|\mathcal{F}_{t,0}] \\ &\leq \mathbb{E}[G_{t,0}|\mathcal{F}_{t,0}] - \mathbb{E}[G_{t,k_{\text{in}}}| \mathcal{F}_{t,0}] \\ &\quad + \frac{2}{\alpha}L^2\gamma^2\left(1 + \frac{\alpha^2k_{\text{in}}}{\mathbf{b}}\right)\sum_{k=1}^{k_{\text{in}}}\mathbb{E}[\|H_{t,k}\|^2|\mathcal{F}_{t,0}]. \end{aligned}$$

Summing from $t = 1$ to $t = k_{\text{out}}$, dividing by $k_{\text{out}}k_{\text{in}}$, and taking expectation we get:

$$\begin{aligned} \frac{1}{k_{\text{out}}k_{\text{in}}}\sum_{t=1}^{k_{\text{out}}}\sum_{k=0}^{k_{\text{in}}-1}\mathbb{E}[G_{t,k}] &\leq \frac{2}{\alpha k_{\text{out}}k_{\text{in}}}\mathbb{E}[G_{1,0}] \\ &\quad + \frac{4}{\alpha^2k_{\text{out}}k_{\text{in}}}L^2\gamma^2\left(1 + \frac{\alpha^2k_{\text{in}}}{\mathbf{b}}\right)\sum_{t=1}^{k_{\text{out}}}\sum_{k=1}^{k_{\text{in}}}\mathbb{E}[\|H_{t,k}\|^2]. \end{aligned}$$

We used that $G_{t,k_{\text{in}}} = G_{t+1,0}$. By denoting (τ, K) a uniform random variable on $[k_{\text{out}}]^* \times [k_{\text{in}} - 1]$ – independent of the path $\{\hat{S}_{t,k}, t \in [k_{\text{out}}]^*, k \in [k_{\text{in}}]\}$, we have

$$\mathbb{E}[G_{\tau,K}] \leq \frac{2}{\alpha k_{\text{out}}k_{\text{in}}}\mathbb{E}[G_{1,0}] + \frac{4}{\alpha^2}L^2\gamma^2\left(1 + \frac{\alpha^2k_{\text{in}}}{\mathbf{b}}\right)\mathbb{E}[\|H_{\tau,K+1}\|^2].$$

From Theorem 4, this yields (note that $\alpha = (1 + \omega)^{-1}$ and $\mathbf{b} \geq k_{\text{in}}/(1 + \omega)^2$)

$$\mathbb{E}[G_{\tau,K}] \leq \frac{2(1 + \omega)}{k_{\text{out}}k_{\text{in}}}\mathbb{E}[G_{1,0}] + \gamma\frac{16(1 + \omega)^2L^2}{v_{\text{min}}k_{\text{in}}k_{\text{out}}}\left(\mathbb{W}(\hat{S}_{1,0}) - \min \mathbb{W}\right).$$

□

10 Supplement to the numerical section

This section gathers additional details concerning the models used in our numerical experiments. Namely, Section 10.1 presents the full derivations for the FedEM algorithm for finite Gaussian Mixture Models, and Section 10.2 provides the detailed pseudo-code for the FedMissEM algorithm for federated missing values imputation introduced in Section 4 and provides the necessary information to request access to the data we used on the eBird platform [1].

10.1 Gaussian Mixture Model

Let y_1, \dots, y_N be N \mathbb{R}^p -valued observations; they are modeled as the realization of a vector (Y_1, \dots, Y_N) with distribution defined as follows:

- conditionally to a $\{1, \dots, L\}$ -valued vector of random variables (Z_1, \dots, Z_N) , (Y_1, \dots, Y_N) are independent; and the conditional distribution of Y_i is $\mathcal{N}_p(\mu_{Z_i}, \Sigma)$.

- the r.v. (Z_1, \dots, Z_n) are i.i.d. with multinomial distribution of size 1 and with probabilities π_1, \dots, π_L .

Equivalently, the random variables (Y_1, \dots, Y_N) are independent with distribution $\sum_{\ell=1}^L \pi_\ell \mathcal{N}_p(\mu_\ell, \Sigma)$. For $1 \leq i \leq N$, the negative log-likelihood is given up to a constant term by

$$\mathcal{L}_i(\theta) = -\langle \mathbf{s}_i(z), \phi(\theta) \rangle + \psi(\theta),$$

where, denoting $\mathbb{1}_{\{l\}}(z)$ the indicator function equal to 1 if $z = l$ and 0 otherwise:

$$\mathbf{s}_i(z) \stackrel{\text{def}}{=} \begin{pmatrix} \mathbb{1}_{\{1\}}(z) \\ \vdots \\ \mathbb{1}_{\{L\}}(z) \\ Y_i \mathbb{1}_{\{1\}}(z) \\ \vdots \\ Y_i \mathbb{1}_{\{L\}}(z) \\ \text{Vec}(Y_i Y_i^\top) \end{pmatrix}, \phi(\theta) \stackrel{\text{def}}{=} \begin{pmatrix} \log(\pi_1) - \mu_1^\top \Sigma^{-1} \mu_1 \\ \vdots \\ \log(\pi_L) - \mu_L^\top \Sigma^{-1} \mu_L \\ \Sigma^{-1} \mu_1 \\ \vdots \\ \Sigma^{-1} \mu_L \\ -\frac{1}{2} \text{Vec}(\Sigma^{-1}) \end{pmatrix}, \psi(\theta) \stackrel{\text{def}}{=} 0. \quad (58)$$

In (58), for $M \in \mathbb{R}^{p \times p}$, $\text{Vec}(M)$ is the column-wise vectorization of M . The goal is to estimate the parameter $\theta = (\pi_1, \dots, \pi_L, \mu_1, \dots, \mu_L, \Sigma)$ by maximizing the normalized negative log-likelihood:

$$\min_{\theta} F(\theta) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_i(\theta). \quad (59)$$

Classical EM algorithm We use the EM algorithm: in the Expectation (E) step, using the current value of the iterate θ_{curr} , we compute a majorizing function $\theta \mapsto Q(\theta, \theta_{\text{curr}})$ given up to an additive constant by

$$Q(\theta, \theta_{\text{curr}}) = -\langle \bar{\mathbf{s}}(\theta_{\text{curr}}), \phi(\theta) \rangle + \psi(\theta),$$

where $\bar{\mathbf{s}}(\theta_{\text{curr}}) = \frac{1}{n} \sum_{i=1}^n \bar{\mathbf{s}}_i(\theta)$, and for any $i \in [n]^*$, $\bar{\mathbf{s}}_i(\theta)$ is the conditional expectation of the complete data sufficient statistics:

$$\bar{\mathbf{s}}_i(\theta) = \begin{pmatrix} \bar{\rho}_{i,1}(\theta) \\ \vdots \\ \bar{\rho}_{i,L}(\theta) \\ Y_i \bar{\rho}_{i,1}(\theta) \\ \vdots \\ Y_i \bar{\rho}_{i,L}(\theta) \\ \text{Vec}(Y_i Y_i^\top) \end{pmatrix}, \text{ where for } \ell \in [L]^*, \bar{\rho}_{i,\ell}(\theta) \stackrel{\text{def}}{=} \frac{\pi_\ell \mathcal{N}_p(\mu_\ell, \Sigma)[Y_i]}{\sum_{u=1}^L \pi_u \mathcal{N}_p(\mu_u, \Sigma)[Y_i]}. \quad (60)$$

In (60), $\mathcal{N}_p(\mu, \Sigma)[y]$ is the density function of the distribution $\mathcal{N}_p(\mu, \Sigma)$ evaluated at y .

In the optimization step (M-step), a new value of θ_{curr} is computed as a minimizer of $\theta \mapsto Q(\theta, \theta_{\text{curr}})$. Let $\mathbf{s} = (s^{(1)}, s^{(2)}, s^{(3)}) \in \mathbb{R}^L \times \mathbb{R}^{pL} \times \mathbb{R}^{p^2}$; we write $\langle \mathbf{s}, \phi(\theta) \rangle = \sum_{j=1}^3 \langle s^{(j)}, \phi^{(j)}(\theta) \rangle$ where the functions $\phi^{(j)}$ are defined by

$$\phi^{(1)}(\theta) \stackrel{\text{def}}{=} \begin{pmatrix} \log(\pi_1) - \mu_1^\top \Sigma^{-1} \mu_1 \\ \vdots \\ \log(\pi_L) - \mu_L^\top \Sigma^{-1} \mu_L \end{pmatrix}, \phi^{(2)}(\theta) \stackrel{\text{def}}{=} \begin{pmatrix} \Sigma^{-1} \mu_1 \\ \vdots \\ \Sigma^{-1} \mu_L \end{pmatrix}, \text{ and } \phi^{(3)}(\theta) \stackrel{\text{def}}{=} \frac{1}{2} \text{Vec}(\Sigma^{-1}). \quad (61)$$

Remember that $\mathbb{T}(\mathbf{s}) = \text{argmin}_{\theta \in \Theta} -\langle \mathbf{s}, \phi(\theta) \rangle + \psi(\theta)$. We obtain $\mathbb{T}(\mathbf{s}) = \{\pi_\ell, \mu_\ell, \ell = 1, \dots, L; \Sigma\}$ with

$$\pi_\ell \stackrel{\text{def}}{=} \frac{s^{(1),\ell}}{\sum_{u=1}^L s^{(1),u}}, \quad (62)$$

$$\mu_\ell \stackrel{\text{def}}{=} \frac{s^{(2),\ell}}{s^{(1),\ell}}, \quad (63)$$

$$\Sigma \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n Y_i Y_i^\top - \sum_{\ell=1}^L s^{(1),\ell} \mu_\ell \mu_\ell^\top. \quad (64)$$

The expressions of π_ℓ, μ_ℓ are easily obtained; we provide details for the covariance matrix. We have for any symmetric matrix H

$$\begin{aligned} \ln \frac{\det(\Gamma + H)}{\det(\Gamma)} &= \ln \det(I + \Gamma^{-1}H) = \ln(1 + \text{Tr}(\Gamma^{-1}H) + o(\|H\|)) \\ &= \text{Tr}(\Gamma^{-1}H) + o(\|H\|) = \langle H, \Gamma^{-1} \rangle + o(\|H\|) \end{aligned}$$

$T(s)$ depends on Γ through the function

$$G(\Gamma) \stackrel{\text{def}}{=} -\frac{1}{2} \ln \det(\Gamma) + \frac{1}{2} \left\langle \Gamma, \frac{1}{n} \sum_{i=1}^n Y_i Y_i^\top + \sum_{\ell=1}^L s^{(1),\ell} \mu_\ell \mu_\ell^\top \right\rangle - \left\langle \Gamma, \sum_{\ell=1}^L \mu_\ell (s^{(2),\ell})^\top \right\rangle.$$

Therefore

$$\Sigma = \Gamma^{-1} \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n Y_i Y_i^\top - \sum_{\ell=1}^L s^{(1),\ell} \mu_\ell \mu_\ell^\top$$

since $\mu_\ell = s^{(2),\ell} / s^{(1),\ell}$.

Algorithm 5: Classical EM algorithm for mixture of Gaussians

- 1: **Input:** $k_{\max} \in \mathbb{N}, X, \hat{S}_0, \hat{\theta}_0$
 - 2: **Output:** The sequence of statistics: $\{\hat{S}_k, k \in [k_{\max}]\}$; the sequence of parameters $\{\hat{\theta}_k, k \in [k_{\max}]\}$
 - 3: **for** $k = 0, \dots, k_{\max} - 1$ **do**
 - 4: *Expectation step:* compute conditional expectations given current parameter $\hat{\theta}_k$: Set $\hat{S}_{k+1} = \frac{1}{N} \sum_{i=1}^N \bar{s}_i(\hat{\theta}^k)$
 - 5: *Maximization step:* update parameter $\hat{\theta}_{k+1}$ based on current statistics \hat{S}_{k+1} according to update rule (62)
 - 6: **end for**
-

In the federated setting, the data is distributed across n local servers. For all $c \in [n]^*$, the c -th server possesses an N_c sample (Y_1, \dots, Y_{N_c}) , and $N_c \geq 1$ and $\sum_{c=1}^n N_c = N$. The pseudo code for the FedEM algorithm is given in Algorithm 6.

10.2 Federated missing values imputation

FedMissEM algorithm. Algorithm 7 provides the pseudo-code for the Federated EM algorithm for missing values imputation.

eBird data information. In our experiments, we used a sample of the eBird data set [1], provided upon request by the Cornell Lab of Ornithology. We are not allowed to disclose the data itself, but we provide here the details to reproduce our experiments on the same data set, after requesting access on the eBird platform (<https://ebird.org/data/request>). We selected the counts recorded anywhere in France, between January 2000 and September 2020, for two different species: the Mallard and the Common Buzzard. These two species were analyzed independently (see Section 4); the corresponding code is also available as supplementary material.

Algorithm 6: Federated EM algorithm for distributed GMM without compression

- 1: **Input:** $k_{\max} \in \mathbb{N}$; for $c \in [n]^*$, $V_{0,c} \in \mathbb{R}^{L+pL}$; $\widehat{S}_0 \in \mathbb{R}^{L+pL}$; $\widehat{\theta}_0 \in \mathbb{R}^L \times (\mathbb{R}^p)^L \times \mathbb{R}^{p \times p}$; a positive sequence $\{\gamma_{k+1}, k \in [k_{\max} - 1]\}$; α
 - 2: **Output:** The FedEMsequence: $\{\widehat{S}_k, k \in [k_{\max}]\}$
 - 3: **for** $k = 0, \dots, k_{\max} - 1$ **do**
 - 4: **for** $c = 1, \dots, n$ **do**
 - 5: (agent # i , locally)
 - 6: Sample a batch $\mathcal{I}_{k,c} \subset [N_c]$
 - 7: Set $S_{k+1,c} = \frac{1}{|\mathcal{I}_{k,c}|} \sum_{i \in \mathcal{I}_{k,c}} \bar{s}_i(\widehat{\theta}_k)$, where \bar{s}_i is defined in (60)
 - 8: Set $\Delta_{k+1,c} = S_{k+1,c} - \widehat{S}_k - V_{k,c}$
 - 9: Update $V_{k+1,c} = V_{k,c} + \alpha \text{Quant}(\Delta_{k+1,c})$
 - 10: Send $\text{Quant}(\Delta_{k+1,c})$ to the controller
 - 11: **end for**
 - 12: (the controller)
 - 13: Compute $H_{k+1} = V_k + \frac{1}{n} \sum_{c=1}^n \text{Quant}(\Delta_{k+1,c})$
 - 14: Set $\widehat{S}_{k+1} = \widehat{S}_k + \gamma_{k+1} H_{k+1}$
 - 15: Set $V_{k+1} = V_k + \alpha n^{-1} \sum_{c=1}^n \text{Quant}(\Delta_{k+1,c})$.
 - 16: Send \widehat{S}_{k+1} and $\widehat{\theta}_{k+1} = \mathbb{T}(\widehat{S}_{k+1})$ to the agents, where $\mathbb{T}(\widehat{S}_{k+1})$ is given by the update rule (62)
 - 17: **end for**
-

Algorithm 7: Federated EM algorithm for distributed missing data imputation

- 1: **Input:** $k_{\max} \in \mathbb{N}$; for $c \in [n]^*$, $V_0^c \in \mathbb{R}^{I \times J}$; $\widehat{S}_0 \in \mathbb{R}^{I \times J}$; a positive sequence $\{\gamma_{k+1}, k \in [k_{\max} - 1]\}$; α ; the quantization function Quant
 - 2: **Output:** The FedEM sequence: $\{\widehat{S}_k, k \in [k_{\max}]\}$
 - 3: **for** $k = 0, \dots, k_{\max} - 1$ **do**
 - 4: **for** $c = 1, \dots, n$ **do**
 - 5: (agent # i , locally)
 - 6: Initialize $S_{k+1,c} = 0$ and $\Delta_{k+1,c} = 0$ everywhere.
 - 7: Sample a minibatch $(\mathcal{I}_k^c, \mathcal{J}_k^c) \subset [I]^* \times [J]^*$
 - 8: **for** $i \in \mathcal{I}_k^c$ **do**
 - 9: **for** $j \in \mathcal{J}_k^c$ **do**
 - 10: Set $(S_{k+1}^c)_{i,j} = \mathbb{1}_{i,j \in \Omega^c} Y_{i,j}^c + (1 - \mathbb{1}_{i,j \in \Omega^c})(\widehat{\theta}_k)_{i,j}$
 - 11: Set $(\Delta_{k+1}^c)_{i,j} = (S_{k+1}^c)_{i,j} - \widehat{S}_{i,j} - (V_k^c)_{i,j}$
 - 12: **end for**
 - 13: **end for**
 - 14: Update $V_{k+1}^c = V_k^c + \alpha \text{Quant}(\Delta_{k+1,c})$
 - 15: Send $\text{Quant}(\Delta_{k+1}^c)$ to the controller
 - 16: **end for**
 - 17: (the controller)
 - 18: Compute $H_{k+1} = V_k + n^{-1} \sum_{c=1}^n \text{Quant}(\Delta_{k+1}^c)$
 - 19: Set $\widehat{S}_{k+1} = \widehat{S}_k + \gamma_{k+1} H_{k+1}$
 - 20: Set $V_{k+1} = V_k + \alpha n^{-1} \sum_{c=1}^n \text{Quant}(\Delta_{k+1}^c)$.
 - 21: Send \widehat{S}_{k+1} and $\widehat{\theta}_{k+1} = \mathbb{T}(\widehat{S}_{k+1})$ to the agents
 - 22: (Note: thresholded SVD for Gaussian model or computed iteratively for a general exponential family model)
 - 23: **end for**
-