



Bidirectional compression for federated learning in heterogeneous setting

Constantin Philippenko
CMAP, École Polytechnique, Institut Polytechnique de Paris
Accenture Labs, Sophia Antipolis

Reviewers:
Mikael Johansson
Jérôme Malick

Ph.D. supervisors: [Aymeric Dieuleveut](#) and [Eric Moulines](#)

Ph.D. defense, September 18th, 2023

Examiners:
Manon Costa
Robert Gower
Martin Jaggi
Kevin Scaman

General introduction

Framework for bidirectional compression

Contributions

I. *Artemis* and the memory mechanism

II. MCM and the preserved update equation

III. Beyond worst-case analysis

Conclusion

General introduction

← Identification - Résultats
Plantes utiles



Ligularia dentata (A.Gray) Hara

Ligulaire dentée

Asteraceae

✓ Valider

 85%

Figure 1: Automatic plant identification from photos using the mobile app [PI@ntNet].

← Identification - Résultats ▾
Plantes utiles



Ligularia dentata (A.Gray) Hara

Ligulaire dentée

Asteraceae

✓ Valider

85%

Goal of machine learning:

Find a mathematical relationship between the input (here the images) and the output (here the name of the plant).

Figure 1: Automatic plant identification from photos using the mobile app [PI@ntNet].

← Identification - Résultats ▾
Plantes utiles



Ligularia dentata (A.Gray) Hara

Ligulaire dentée

Asteraceae

✓ Valider

85%

Goal of machine learning:

Find a mathematical relationship between the input (here the images) and the output (here the name of the plant).

Paradigm of my thesis: data is not centralized on a single location.

Figure 1: Automatic plant identification from photos using the mobile app [PI@ntNet].

← Identification - Résultats ▾
Plantes utiles



Ligularia dentata (A.Gray) Hara

Ligulaire dentée

Asteraceae

✓ Valider

85%

Goal of machine learning:

Find a mathematical relationship between the input (here the images) and the output (here the name of the plant).

Paradigm of my thesis: data is not centralized on a single location.

Privacy

Communication cost

Data heterogeneity

Figure 1: Automatic plant identification from photos using the mobile app [PI@ntNet].

← Identification - Résultats ▾
Plantes utiles



Ligularia dentata (A.Gray) Hara

Ligulaire dentée

Asteraceae

✓ Valider

85%

Goal of machine learning:

Find a mathematical relationship between the input (here the images) and the output (here the name of the plant).

Paradigm of my thesis: data is not centralized on a single location.

Privacy

Communication cost

Data heterogeneity

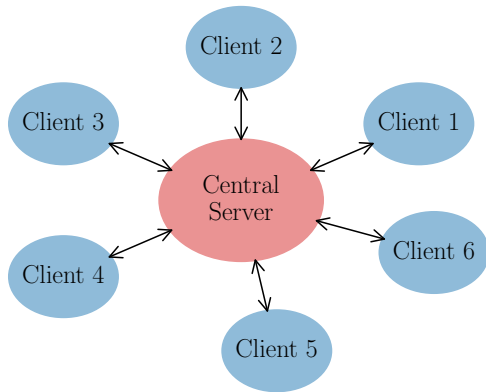
Goal of my thesis:

Focus simultaneously on two challenges: **reducing the cost of communication** and considering a **heterogeneous setting**.

Figure 1: Automatic plant identification from photos using the mobile app [PI@ntNet].

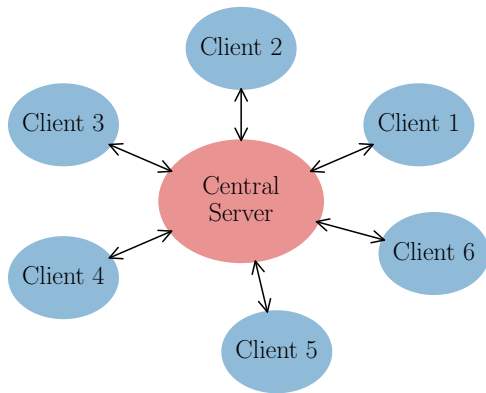
Setting of federated learning:

A central server orchestrate the training.



Setting of federated learning:

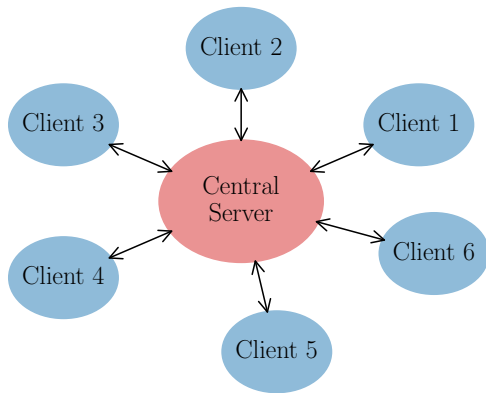
A central server orchestrate the training.



Each client $i \in \mathbb{N}^*$ have access to a “objective function” F_i measuring the error of prediction for a model $w \in \mathbb{R}^d$.

Setting of federated learning:

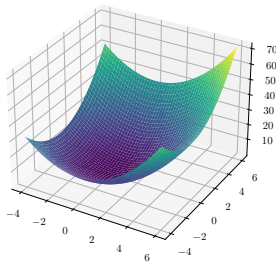
A central server orchestrate the training.



We need to find the optimal model w_* such that:

$$w_* = \operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N F_i(w).$$

$$F_1 : x, y \mapsto x^2 + y^2$$



$$F_2 : x, y \mapsto (1 - \sin(x))^2 + \cos(y)$$

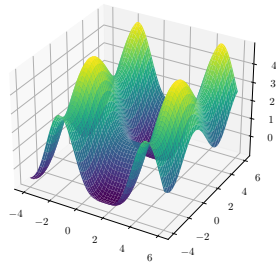
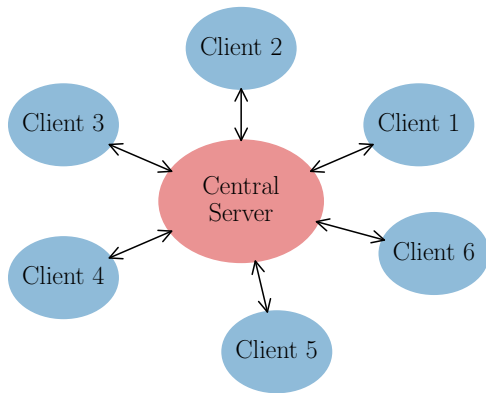


Figure 2: Examples of two objective functions

Each client $i \in \mathbb{N}^*$ have access to a “objective function” F_i measuring the error of prediction for a model $w \in \mathbb{R}^d$.

Setting of federated learning:

A central server orchestrate the training.



Each client $i \in \mathbb{N}^*$ have access to a “objective function” F_i measuring the error of prediction for a model $w \in \mathbb{R}^d$.

We need to find the optimal model w_* such that:

$$w_* = \operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N F_i(w).$$

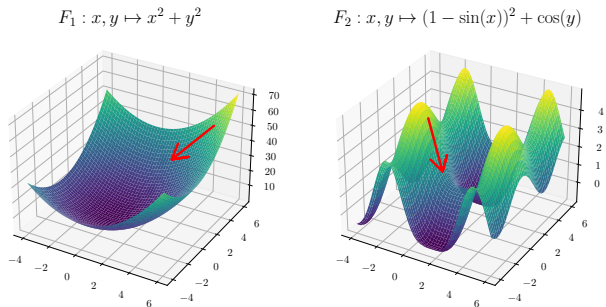


Figure 2: Examples of two objective functions

To find the optimal model w_* , we follow the slope (gradient descent).

Framework for bidirectional compression

Two challenges of Federated Learning



Goal : learning from a set of N clients [MMR⁺17]

$$\min_{w \in \mathbb{R}^d} \left\{ F(w) := \frac{1}{N} \sum_{i=1}^N \underbrace{\mathbb{E}_{z \sim \mathcal{D}_i} [\ell(z, w)]}_{F_i(w)} \right\}.$$

F : global cost function
 F_i : local loss
 N : clients
 d : dimension
 w : model
 \mathcal{D}_i : local data distribution

Global loss

$$F(w) := \frac{1}{N} \sum_{i=1}^N F_i(w)$$

Local loss



Distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N g_k^i(w_{k-1}) \right)$.

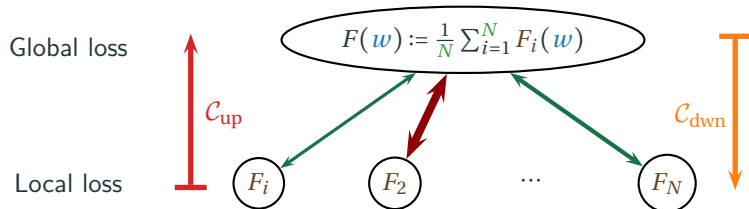
Two challenges of Federated Learning



Goal : learning from a set of N clients [MMR⁺17]

$$\min_{w \in \mathbb{R}^d} \left\{ F(w) := \frac{1}{N} \sum_{i=1}^N \underbrace{\mathbb{E}_{z \sim \mathcal{D}_i} [\ell(z, w)]}_{F_i(w)} \right\}.$$

F : global cost function
 F_i : local loss
 N : clients
 d : dimension
 w : model
 \mathcal{D}_i : local data distribution



→ **Challenge 1:**
reduce communication costs

Distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N g_k^i(w_{k-1}) \right)$.

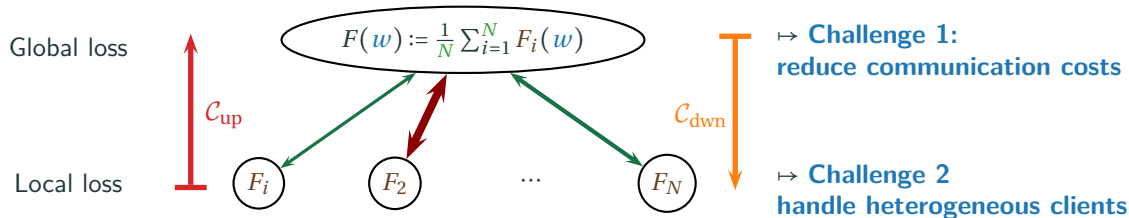
Two challenges of Federated Learning



Goal : learning from a set of N clients [MMR⁺17]

$$\min_{w \in \mathbb{R}^d} \left\{ F(w) := \frac{1}{N} \sum_{i=1}^N \underbrace{\mathbb{E}_{z \sim \mathcal{D}_i} [\ell(z, w)]}_{F_i(w)} \right\}.$$

F : global cost function
 F_i : local loss
 N : clients
 d : dimension
 w : model
 \mathcal{D}_i : local data distribution



Distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N g_k^i(w_{k-1}) \right)$.

- ↳ To limit the number of bits exchanged, we **compress** each signal before transmitting it.
- ↳ **Focus on bidirectional compression** [LLTY20, PD20, TYL⁺19, ZHK19, PD21].

↳ To limit the number of bits exchanged, we **compress** each signal before transmitting it.

↳ **Focus on bidirectional compression** [LLTY20, PD20, TYL⁺19, ZHK19, PD21].

↳ We introduce two compression operators $\mathcal{C}_{\text{dwn}} \downarrow$ and $\mathcal{C}_{\text{up}} \uparrow$.

Compressed distributed SGD:

$$\forall k \in \mathbb{N}, w_{k+1} = w_k - \gamma \mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_{k+1}^i(w_k)) \right).$$

↪ To limit the number of bits exchanged, we **compress** each signal before transmitting it.

↪ **Focus on bidirectional compression** [LLTY20, PD20, TYL⁺19, ZHK19, PD21].

↪ We introduce two compression operators \mathcal{C}_{dwn} \downarrow and \mathcal{C}_{up} \uparrow .

Compressed distributed SGD:

$$\forall k \in \mathbb{N}, w_{k+1} = w_k - \gamma \mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_{k+1}^i(w_k)) \right).$$

Assumption 1 (One assumption to rule them all)

For $\text{dir} \in \{\text{up}, \text{dwn}\}$, there exists a constant $\omega_{\text{dir}} \in \mathbb{R}_+^*$ s.t. \mathcal{C}_{dir} satisfies, for all z in \mathbb{R}^d :

$$\mathbb{E}[\mathcal{C}_{\text{dir}}(z)] = z \quad \text{and} \quad \mathbb{E}[\|\mathcal{C}_{\text{dir}}(z) - z\|^2] \leq \omega_{\text{dir}} \|z\|^2.$$

The compressors are said to be *Unbiased with a Relatively Bounded Variance* (URBV).

1. Sparsification based:
 - Rand-k: keeps k coordinates,
 - p -Sparsification: keeps each coordinate with probability p ,
 - p -partial participation: sends the complete vector with probability p ,
 - Sketching: using a random projection matrix into a lower-dimension space.
2. Quantization based on a codebook:
 - (Stabilized) scalar quantization (coordinate compressed independently),
 - Delaunay quantization.

Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \gamma \mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}} (g_k^i(w_{k-1})) \right)$.

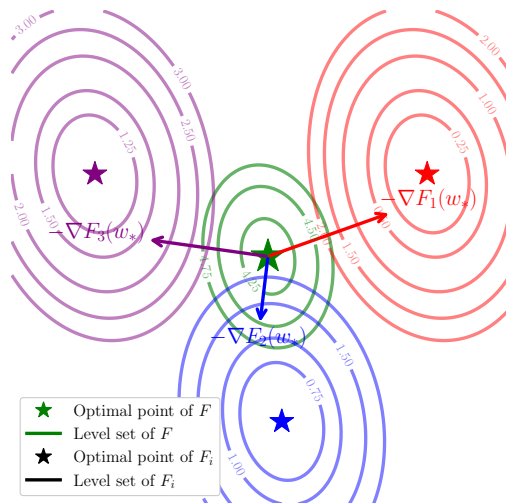
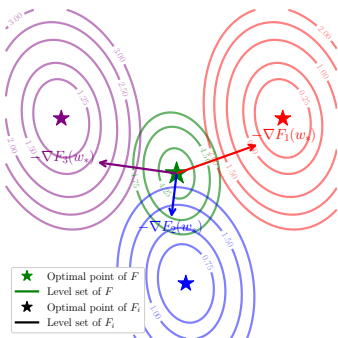


Figure 3: Illustration of heterogeneity on three clients, the objective functions are quadratic. We represent the optimal points, the level set, and the opposite gradient at the optimal point.

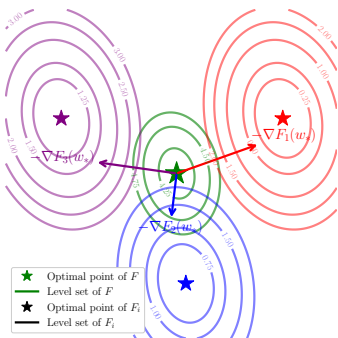
Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(w_{k-1})) \right)$.



Assumption 2 (Bounded gradient at w_*)

There exists an optimal parameter w_* minimizing F (not necessarily unique) and a constant $B \in \mathbb{R}_+$, such that $\frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w_*)\|^2 = B^2$.

Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \gamma \mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(w_{k-1})) \right)$.



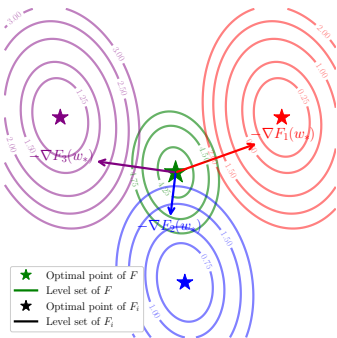
Assumption 2 (Bounded gradient at w_*)

There exists an optimal parameter w_* minimizing F (not necessarily unique) and a constant $B \in \mathbb{R}_+$, such that $\frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w_*)\|^2 = B^2$.

Assumption 3 (Noise over stochastic gradients computation)

The noise over stochastic gradients is zero-centered and its variance is uniformly bounded by a constant $\sigma \in \mathbb{R}_+$, such that for all k in \mathbb{N} , for all z in \mathbb{R}^d we have: $\mathbb{E}[\|g_k(z) - \nabla F(z)\|^2] \leq \sigma^2$.

Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \gamma \mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(w_{k-1})) \right)$.



Assumption 2 (Bounded gradient at w_*)

There exists an optimal parameter w_* minimizing F (not necessarily unique) and a constant $B \in \mathbb{R}_+$, such that $\frac{1}{N} \sum_{i=1}^N \|\nabla F_i(w_*)\|^2 = B^2$.

Assumption 3 (Noise over stochastic gradients computation)

The noise over stochastic gradients is zero-centered and its variance is uniformly bounded by a constant $\sigma \in \mathbb{R}_+$, such that for all k in \mathbb{N} , for all z in \mathbb{R}^d we have: $\mathbb{E}[\|g_k(z) - \nabla F(z)\|^2] \leq \sigma^2$.

Theorem 1 (Convergence of compressed distributed SGD)

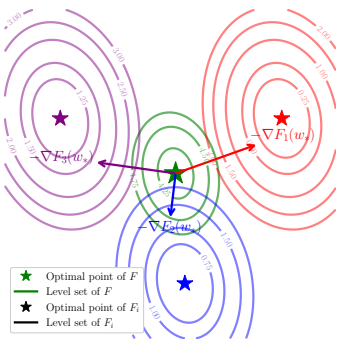
Under A1, A2, A3, if all $(F_i)_{i=1}^N$ are L -smooth, $\mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(w_k)) \right)$ is an unbiased stochastic oracle of $\nabla F(w_{k-1})$ with variance bounded by:

$$\frac{2(\omega_{\text{dwn}} + 1)(\omega_{\text{up}} + 1)\sigma^2}{N} + \frac{4\omega_{\text{dwn}}\omega_{\text{up}}B^2}{N} + 2L\omega_{\text{dwn}}\|w_k - w_*\|^2 \left(1 + \frac{2}{N}\right).$$

From a first theorem to a glance at contributions



Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \gamma \mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(w_{k-1})) \right)$.



Relax the uniform bound

Remove the B^2 -dependence

Theorem 1 (Convergence of compressed distributed SGD)

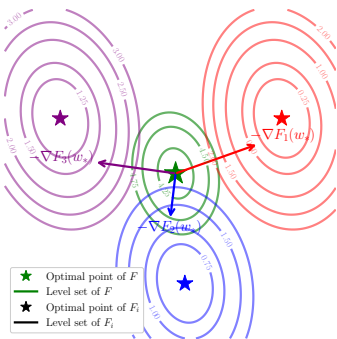
Under A1, A2, A3, if all $(F_i)_{i=1}^N$ are L -smooth, $\mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(w_k)) \right)$ is an unbiased stochastic oracle of $\nabla F(w_{k-1})$ with variance bounded by:

$$\frac{2(\omega_{\text{dwn}} + 1)(\omega_{\text{up}} + 1)\sigma^2}{N} + \frac{4\omega_{\text{dwn}}\omega_{\text{up}}B^2}{N} + 2L\omega_{\text{dwn}}\|w_k - w_*\|^2\left(1 + \frac{2}{N}\right).$$

From a first theorem to a glance at contributions



Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \gamma \mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(w_{k-1})) \right)$.



Remove the w_{dwn} -dependence in the dominant term

Theorem 1 (Convergence of compressed distributed SGD)

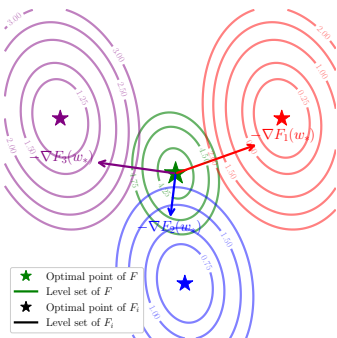
Under A1, A2, A3, if all $(F_i)_{i=1}^N$ are L -smooth, $\mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(w_k)) \right)$ is an unbiased stochastic oracle of $\nabla F(w_{k-1})$ with variance bounded by:

$$\frac{2(\omega_{\text{dwn}} + 1)(\omega_{\text{up}} + 1)\sigma^2}{N} + \frac{4\omega_{\text{dwn}}\omega_{\text{up}}B^2}{N} + 2L\omega_{\text{dwn}}\|w_k - w_*\|^2 \left(1 + \frac{2}{N}\right).$$

From a first theorem to a glance at contributions



Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \gamma \mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(w_{k-1})) \right)$.



Going beyond the worst-case assumption

Theorem 1 (Convergence of compressed distributed SGD)

Under A1, A2, A3, if all $(F_i)_{i=1}^N$ are L -smooth, $\mathcal{C}_{\text{dwn}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(w_k)) \right)$ is an unbiased stochastic oracle of $\nabla F(w_{k-1})$ with variance bounded by:

$$\frac{2(\omega_{\text{dwn}} + 1)(\omega_{\text{up}} + 1)\sigma^2}{N} + \frac{4\omega_{\text{dwn}}\omega_{\text{up}}B^2}{N} + 2L\omega_{\text{dwn}}\|w_k - w_*\|^2\left(1 + \frac{2}{N}\right).$$

Contributions

- I. *Artemis: tight convergence guarantees for bidirectional compression with heterogeneous clients*, P and Dieuleveut, under review at *Journal of Parallel and Distributed Computing*
- II. *MCM: a preserved central model for faster bidirectional compression in distributed settings*, P and Dieuleveut, Neurips 2021
- III. *Convergence rates for distributed, compressed and averaged least-squares regression: application to Federated Learning*, P and Dieuleveut, under review at *Journal of Machine Learning Research*

- I. *Artemis: tight convergence guarantees for bidirectional compression with heterogeneous clients*, P and Dieuleveut, under review at *Journal of Parallel and Distributed Computing*
- II. *MCM: a preserved central model for faster bidirectional compression in distributed settings*, P and Dieuleveut, Neurips 2021
- III. *Convergence rates for distributed, compressed and averaged least-squares regression: application to Federated Learning*, P and Dieuleveut, under review at *Journal of Machine Learning Research*

Table 1: Summary of contributions.

	Bi-compr.	Heterogeneity	LSR	
I.	✓	✓		Interaction between compression and heterogeneity
II.	✓		(✓)	Asympt. cancels impact of down compression
III.		(✓)	✓	Beyond worst-case analysis

- I. *Artemis: tight convergence guarantees for bidirectional compression with heterogeneous clients*, P and Dieuleveut, under review at *Journal of Parallel and Distributed Computing*
- II. *MCM: a preserved central model for faster bidirectional compression in distributed settings*, P and Dieuleveut, Neurips 2021
- III. *Convergence rates for distributed, compressed and averaged least-squares regression: application to Federated Learning*, P and Dieuleveut, under review at *Journal of Machine Learning Research*

Table 1: Summary of contributions.

	Bi-compr.	Heterogeneity	LSR	
I.	✓	✓		Interaction between compression and heterogeneity
II.	✓		(✓)	Asympt. cancels impact of down compression
III.		(✓)	✓	Beyond worst-case analysis

Not included in my manuscript: *FLamby: Datasets and benchmarks for cross-silo federated learning in realistic healthcare settings*, Ogier du Terrail, [...] P, [...] Andreux, Neurips 2022.

I. Artemis and the memory mechanism

We make standard assumptions on $F: \mathbb{R}^d \rightarrow \mathbb{R}$.

Assumption 4 (Cocoercivity)

All $(g_k^i)_{i=1}^N$ stochastic gradient are L -cocoercive in quadratic mean.

We make standard assumptions on $F: \mathbb{R}^d \rightarrow \mathbb{R}$.

Assumption 4 (Cocoercivity)

All $(g_k^i)_{i=1}^N$ stochastic gradient are L -cocoercive in quadratic mean.

Assumption 5 (Strong-convexity)

F is strongly-convex.

We make standard assumptions on $F: \mathbb{R}^d \rightarrow \mathbb{R}$.

Assumption 4 (Cocoercivity)

All $(g_k^i)_{i=1}^N$ stochastic gradient are L -cocoercive in quadratic mean.

Assumption 5 (Strong-convexity)

F is strongly-convex.

Extension: We extend our results to the **convex case**.

We make standard assumptions on $F: \mathbb{R}^d \rightarrow \mathbb{R}$.

Assumption 4 (Cocoercivity)

All $(g_k^i)_{i=1}^N$ stochastic gradient are L -cocoercive in quadratic mean.

Assumption 5 (Strong-convexity)

F is strongly-convex.

Extension: We extend our results to the **convex case**.

Assumption 6 (Noise over stochastic gradients computation)

The noise over stochastic gradients for a mini-batch of size b , is bounded at w_* :

$$\exists \sigma_* \in \mathbb{R}_+, \quad \forall k \in \mathbb{N}, \quad \forall i \in \llbracket 1, N \rrbracket, \quad \forall w \in \mathbb{R}^d: \quad E[\|g_k^i(w_*) - \nabla F_i(w_*)\|^2] \leq \sigma_*^2/b.$$

[As in GLQ⁺19, DDB20]

Compressed distributed SGD: $w_k = w_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i) \right)$

Consequence of clients' heterogeneity: $\lim_{k \rightarrow +\infty} g_{k+1}^i(w_*) \neq 0$.

Compressed distributed SGD: $w_k = w_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i) \right)$

Consequence of clients' heterogeneity: $\lim_{k \rightarrow +\infty} g_{k+1}^i(w_*) \neq 0$.

Goal: Compress a quantity that goes to 0

Solution: Compute (on the server and the worker independently) a “memory” h_k^i s.t.

$$h_k^i \xrightarrow[k \rightarrow \infty]{} \nabla F_i(w_*).$$

Compressed distributed SGD: $w_k = w_{k-1} - \gamma C_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i) \right)$

Consequence of clients' heterogeneity: $\lim_{k \rightarrow +\infty} g_{k+1}^i(w_*) \neq 0$.

Goal: Compress a quantity that goes to 0

Solution: Compute (on the server and the worker independently) a “memory” h_k^i s.t.

$$h_k^i \xrightarrow[k \rightarrow \infty]{} \nabla F_i(w_*).$$

⇒ The update equation becomes:

$$w_k = w_{k-1} - \gamma C_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i - h_{k-1}^i) + h_{k-1}^i \right)$$
$$h_k^i = h_{k-1}^i + \alpha C_{\text{up}}(g_k^i - h_{k-1}^i),$$

where α is the memory's learning rate.

Compressed distributed SGD: $w_k = w_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i) \right)$

Consequence of clients' heterogeneity: $\lim_{k \rightarrow +\infty} g_{k+1}^i(w_*) \neq 0$.

Goal: Compress a quantity that goes to 0

Solution: Compute (on the server and the worker independently) a **"memory"** h_k^i s.t.

$$h_k^i \xrightarrow[k \rightarrow \infty]{} \nabla F_i(w_*).$$

⇒ The update equation becomes:

$$w_k = w_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i - h_{k-1}^i) + h_{k-1}^i \right)$$
$$h_k^i = h_{k-1}^i + \alpha \mathcal{C}_{\text{up}}(g_k^i - h_{k-1}^i),$$

where α is the memory's learning rate.

⇒ **Introducing this uplink memory mechanism is crucial to handle data heterogeneity, see Theorem 2.**

Compressed distributed SGD: $w_k = w_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i) \right)$

Consequence of clients' heterogeneity: $\lim_{k \rightarrow +\infty} g_{k+1}^i(w_*) \neq 0$.

Goal: Compress a quantity that goes to 0

Solution: Compute (on the server and the worker independently) a “memory” h_k^i s.t.

$$h_k^i \xrightarrow[k \rightarrow \infty]{} \nabla F_i(w_*).$$

⇒ The update equation becomes:

$$w_k = w_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i - h_{k-1}^i) + h_{k-1}^i \right)$$
$$h_k^i = h_{k-1}^i + \alpha \mathcal{C}_{\text{up}}(g_k^i - h_{k-1}^i),$$

γ : SGD step-size

α : memory's learning rate

where α is the memory's learning rate.

⇒ **Introducing this uplink memory mechanism is crucial to handle data heterogeneity, see Theorem 2.**

Theorem 2 (Convergence of Artemis)

Under A1-2 and A4-6, for a step size γ under some conditions, for a learning rate α and for any k in \mathbb{N} ,

$$\mathbb{E}[\|w_k - w_*\|^2] \leq (1 - \gamma\mu)^k \text{Bias}^2 + 2\gamma \frac{\text{Var}}{\mu N},$$

with:	Variant	Var
$\alpha = 0$		$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)(\sigma_*^2 + B^2)$
$\alpha(\omega_{\text{up}} + 1) = 1/2$		$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)\sigma_*^2$

Theorem 2 (Convergence of Artemis)

Under A1-2 and A4-6, for a step size γ under some conditions, for a learning rate α and for any k in \mathbb{N} ,

$$\mathbb{E}[\|w_k - w_*\|^2] \leq (1 - \gamma\mu)^k \text{Bias}^2 + 2\gamma \frac{\text{Var}}{\mu N},$$

with:	Variat	Var
$\alpha = 0$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)(\sigma_*^2 + B^2)$	
$\alpha(\omega_{\text{up}} + 1) = 1/2$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)\sigma_*^2$	

- Linear rate up to a constant of the order of Var

Theorem 2 (Convergence of Artemis)

Under A1-2 and A4-6, for a step size γ under some conditions, for a learning rate α and for any k in \mathbb{N} ,

$$\mathbb{E}[\|w_k - w_*\|^2] \leq (1 - \gamma\mu)^k \text{Bias}^2 + 2\gamma \frac{\text{Var}}{\mu N},$$

with:

Variant	Var
$\alpha = 0$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)(\sigma_*^2 + B^2)$
$\alpha(\omega_{\text{up}} + 1) = 1/2$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)\sigma_*^2$

- Linear rate up to a constant of the order of Var
- The variance (Var) increases with the compression level.

Theorem 2 (Convergence of Artemis)

Under A1-2 and A4-6, for a step size γ under some conditions, for a learning rate α and for any k in \mathbb{N} ,

$$\mathbb{E}[\|w_k - w_*\|^2] \leq (1 - \gamma\mu)^k \text{Bias}^2 + 2\gamma \frac{\text{Var}}{\mu N},$$

with:

Variant	Var
$\alpha = 0$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)(\sigma_*^2 + B^2)$
$\alpha(\omega_{\text{up}} + 1) = 1/2$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)\sigma_*^2$

- **Linear rate** up to a constant of the order of **Var**
- The variance (**Var**) increases with the compression level.
- When $B^2 \neq 0$ (non-i.i.d. settings), if $\sigma_*^2 = 0$, then using memory ($\alpha \neq 0$) leads to linear convergence

Theorem 2 (Convergence of Artemis)

Under A1-2 and A4-6, for a step size γ under some conditions, for a learning rate α and for any k in \mathbb{N} ,

$$\mathbb{E}[\|w_k - w_*\|^2] \leq (1 - \gamma\mu)^k \text{Bias}^2 + 2\gamma \frac{\text{Var}}{\mu N},$$

with:

Variant	Var
$\alpha = 0$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)(\sigma_*^2 + B^2)$
$\alpha(\omega_{\text{up}} + 1) = 1/2$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)\sigma_*^2$

- Linear rate up to a constant of the order of Var
- The variance (Var) increases with the compression level.
- When $B^2 \neq 0$ (non-i.i.d. settings), if $\sigma_*^2 = 0$, then using memory ($\alpha \neq 0$) leads to linear convergence
- If $B^2 = 0$ (i.i.d. settings), the memory is useless

Theorem 2 (Convergence of Artemis)

Under A1-2 and A4-6, for a step size γ under some conditions, for a learning rate α and for any k in \mathbb{N} ,

$$\mathbb{E}[\|w_k - w_*\|^2] \leq (1 - \gamma\mu)^k \text{Bias}^2 + 2\gamma \frac{\text{Var}}{\mu N},$$

with:

Variant	Var
$\alpha = 0$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)(\sigma_*^2 + B^2)$
$\alpha(\omega_{\text{up}} + 1) = 1/2$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)\sigma_*^2$

- **Linear rate** up to a constant of the order of **Var**
- The variance (**Var**) increases with the compression level.
- When $B^2 \neq 0$ (non-i.i.d. settings), if $\sigma_*^2 = 0$, then using memory ($\alpha \neq 0$) leads to linear convergence
- If $B^2 = 0$ (i.i.d. settings), the memory is useless
- Recovers classical SGD rate in the absence of compression

Theorem 2 (Convergence of Artemis)

Under A1-2 and A4-6, for a step size γ under some conditions, for a learning rate α and for any k in \mathbb{N} ,

$$\mathbb{E}[\|w_k - w_*\|^2] \leq (1 - \gamma\mu)^k \text{Bias}^2 + 2\gamma \frac{\text{Var}}{\mu N},$$

with:

Variant	Var
$\alpha = 0$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)(\sigma_*^2 + B^2)$
$\alpha(\omega_{\text{up}} + 1) = 1/2$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)\sigma_*^2$

- **Linear rate** up to a constant of the order of **Var**
- The variance (**Var**) increases with the compression level.
- When $B^2 \neq 0$ (non-i.i.d. settings), if $\sigma_*^2 = 0$, then using memory ($\alpha \neq 0$) leads to linear convergence
- If $B^2 = 0$ (i.i.d. settings), the memory is useless
- Recovers classical SGD rate in the absence of compression

Theorem 3 (Lower bound on the variance for linear compressors)

Under A1-2 and A4-6, for $\gamma, \alpha_{\text{up}}, E$ given in Theorem 2, for Θ_k the distribution of w_k .

There exists a limit distribution $\pi_{\gamma, \alpha}$ s.t. for any $k \geq 1$, for C_0 a constant:

$$\mathcal{W}_2(\Theta_k, \pi_{\gamma, \alpha}) \leq (1 - \gamma\mu)^k C_0.$$

Furthermore:

$$\mathbb{E}[\|w_k - w_*\|^2] \xrightarrow[k \rightarrow \infty]{} \mathbb{E}_{w \sim \pi_{\gamma, \alpha}}[\|w - w_*\|^2]$$

which is lower bounded s.t.:

$$\mathbb{E}_{w \sim \pi_{\gamma, \alpha}}[\|w - w_*\|^2] \underset{\gamma \rightarrow 0}{=} \Omega(\gamma \text{Var} / \mu N).$$

Theorem 2 (Convergence of Artemis)

Under A1-2 and A4-6, for a step size γ under some conditions, for a learning rate α and for any k in \mathbb{N} ,

$$\mathbb{E}[\|w_k - w_*\|^2] \leq (1 - \gamma\mu)^k \text{Bias}^2 + 2\gamma \frac{\text{Var}}{\mu N},$$

with:

Variant	Var
$\alpha = 0$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)(\sigma_*^2 + B^2)$
$\alpha(\omega_{\text{up}} + 1) = 1/2$	$(\omega_{\text{down}} + 1)(\omega_{\text{up}} + 1)\sigma_*^2$

- **Linear rate** up to a constant of the order of **Var**
- The variance (**Var**) increases with the compression level.
- When $B^2 \neq 0$ (non-i.i.d. settings), if $\sigma_*^2 = 0$, then using memory ($\alpha \neq 0$) leads to linear convergence
- If $B^2 = 0$ (i.i.d. settings), the memory is useless
- Recovers classical SGD rate in the absence of compression

Theorem 3 (Lower bound on the variance for linear compressors)

Under A1-2 and A4-6, for $\gamma, \alpha_{\text{up}}, E$ given in Theorem 2, for Θ_k the distribution of w_k .

There exists a limit distribution $\pi_{\gamma, \alpha}$ s.t. for any $k \geq 1$, for C_0 a constant:

$$\mathcal{W}_2(\Theta_k, \pi_{\gamma, \alpha}) \leq (1 - \gamma\mu)^k C_0.$$

Furthermore:

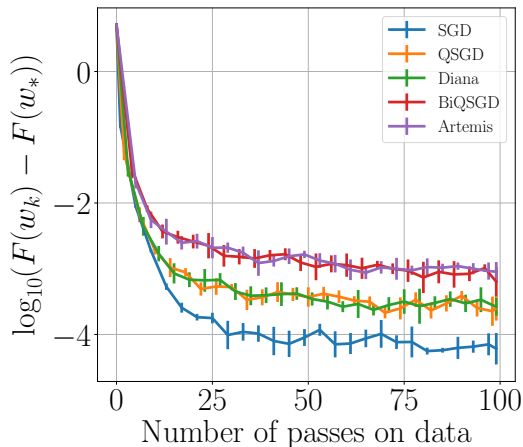
$$\mathbb{E}[\|w_k - w_*\|^2] \xrightarrow[k \rightarrow \infty]{} \mathbb{E}_{w \sim \pi_{\gamma, \alpha}}[\|w - w_*\|^2]$$

which is lower bounded s.t.:

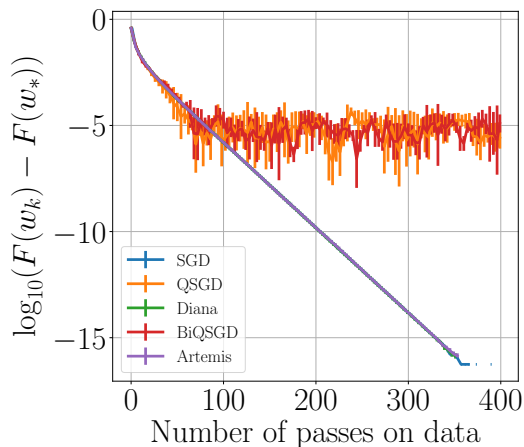
$$\mathbb{E}_{w \sim \pi_{\gamma, \alpha}}[\|w - w_*\|^2] \underset{\gamma \rightarrow 0}{=} \Omega(\gamma \text{Var} / \mu N).$$

The quadratic increase in the variance is not an artifact of the proof!

- Left: illustration of the saturation when $\sigma_*^2 \neq 0$ and data is i.i.d.
- Right: illustration of the memory benefits when $\sigma_*^2 = 0$ but with non-i.i.d. data.



(a) Least-square reg. (i.i.d.): $\sigma_*^2 \neq 0$



(b) Logistic reg. (non-i.i.d.): $\sigma_*^2 = 0$.

Figure 4: Synthetic datasets

- Left: almost homogeneous clients.
- Right: heterogeneous clients.

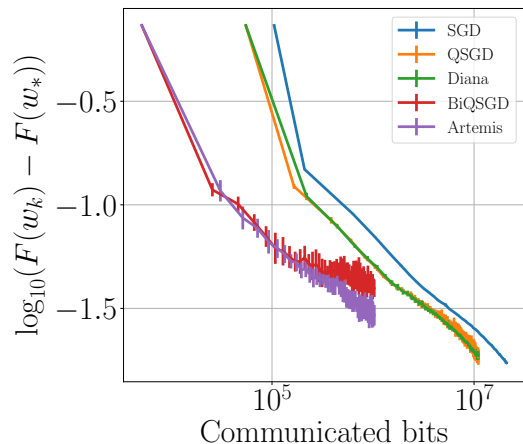


Figure 5: Superconduct (LSR), $b = 64$

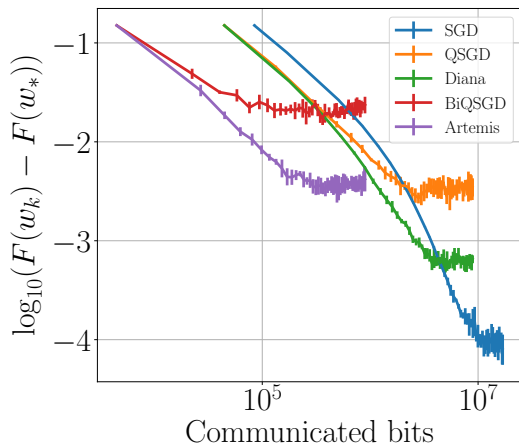


Figure 6: Quantum (LR), $b = 256$

- Stochastic gradient descent: $\sigma_* \neq 0$.

Take-away 1

- ***Bidirectional compression*** to reduce the communication cost.

Take-away 1

- ***Bidirectional compression*** to reduce the communication cost.

Take-away 2

- Primary factor: ***noise σ_* on the gradient*** computed on the ***optimal point***.
- Key impact of ***memory*** on ***non-i.i.d. data***.

Take-away 1

- **Bidirectional compression** to reduce the communication cost.

Take-away 2

- Primary factor: **noise σ_* on the gradient** computed on the **optimal point**.
- Key impact of **memory** on **non-i.i.d. data**.

Take-away 3

- **Lower bound** on the asymptotic variance.

II. MCM and the preserved update equation

Classical approach - degrade the model on the central server.

$$w_k = w_{k-1} - \gamma C_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i(w_{k-1})) \right).$$

The gradient is taken at the point w_k held by the central server
[LLTY20, PD20, TYL⁺19, ZHK19].

Classical approach - **degrade the model on the central server.**

$$w_k = w_{k-1} - \gamma C_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i(w_{k-1})) \right).$$

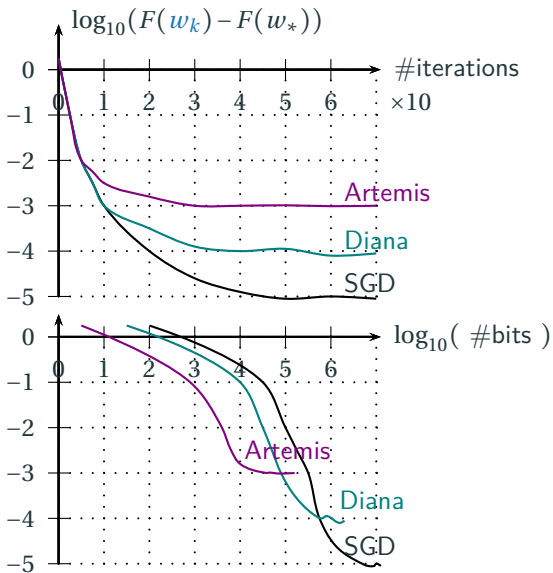
The gradient is taken at the point w_k held by the central server [LLTY20, PD20, TYL⁺19, ZHK19].

New approach - **preserve the model on the central server.**

$$\begin{aligned} w_k &= w_{k-1} - \gamma \frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i(\hat{w}_{k-1})) \\ \hat{w}_k &= \hat{w}_{k-1} - \gamma C_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right). \end{aligned} \tag{1}$$

The gradient is taken at a random point \hat{w}_k s.t. $\mathbb{E}[\hat{w}_k | w_k] = w_k$.

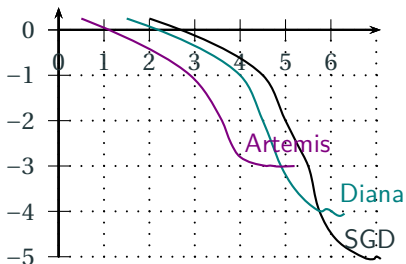
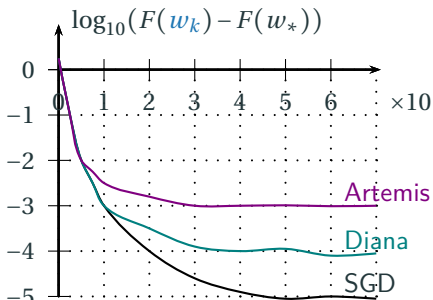
Classical approach



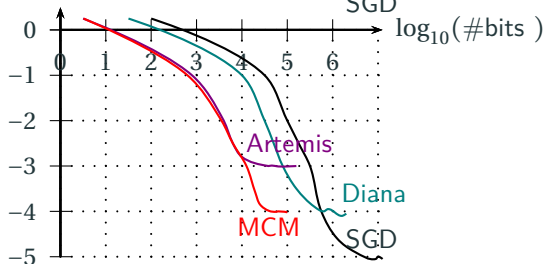
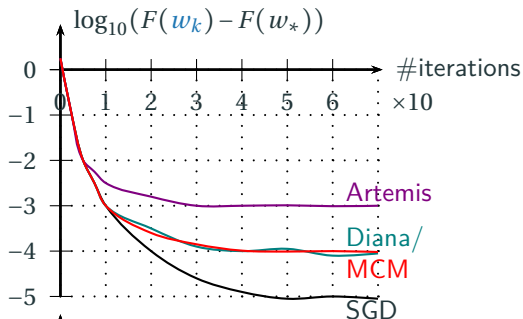
What do we hope for? (using a constant step-size γ)



Classical approach



New approach



We introduce a *downlink memory term* $(H_k)_{k \in \mathbb{N}}$:

1. available on both clients and central server
2. the difference Ω_k between the model and this memory is compressed and exchanged
3. the local model is reconstructed from this information

$$\begin{cases} w_k &= w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right) \\ \Omega_k &= w_k - H_{k-1} \\ \hat{w}_k &= H_{k-1} + \mathcal{C}_{\text{down}}(\Omega_k) \\ H_k &= H_{k-1} + \alpha_{\text{down}} \mathcal{C}_{\text{down}}(\Omega_k). \end{cases} \quad (2)$$

We introduce a *downlink memory term* $(H_k)_{k \in \mathbb{N}}$:

1. available on both clients and central server
2. the difference Ω_k between the model and this memory is compressed and exchanged
3. the local model is reconstructed from this information

$$\begin{cases} w_k &= w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right) \\ \Omega_k &= w_k - H_{k-1} \\ \hat{w}_k &= H_{k-1} + \mathcal{C}_{\text{down}}(\Omega_k) \\ H_k &= H_{k-1} + \alpha_{\text{down}} \mathcal{C}_{\text{down}}(\Omega_k). \end{cases} \quad (2)$$

\implies Introducing this memory mechanism is crucial to control the variance of the local model \hat{w}_k .

We introduce a *downlink memory term* $(H_k)_{k \in \mathbb{N}}$:

1. available on both clients and central server
2. the difference Ω_k between the model and this memory is compressed and exchanged
3. the local model is reconstructed from this information

$$\begin{cases} w_k &= w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right) \\ \Omega_k &= w_k - H_{k-1} \\ \hat{w}_k &= H_{k-1} + \mathcal{C}_{\text{down}}(\Omega_k) \\ H_k &= H_{k-1} + \alpha_{\text{down}} \mathcal{C}_{\text{down}}(\Omega_k). \end{cases} \quad (2)$$

⇒ Introducing this memory mechanism is crucial to control the variance of the local model \hat{w}_k .

⚠ We still use the *uplink memory term* (required to tackle the heterogeneous settings). ⚠

We introduce a *downlink memory term* $(H_k)_{k \in \mathbb{N}}$:

1. available on both clients and central server
2. the difference Ω_k between the model and this memory is compressed and exchanged
3. the local model is reconstructed from this information

$$\begin{cases} w_k &= w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right) \\ \Omega_k &= w_k - H_{k-1} \\ \hat{w}_k &= H_{k-1} + \mathcal{C}_{\text{down}}(\Omega_k) \\ H_k &= H_{k-1} + \alpha_{\text{down}} \mathcal{C}_{\text{down}}(\Omega_k). \end{cases} \quad (2)$$

⇒ Introducing this memory mechanism is crucial to control the variance of the local model \hat{w}_k .

⚠ We still use the *uplink memory term* (required to tackle the heterogeneous settings). ⚠

⇒ This is MCM.

Assumption 7 (**Smoothness and convexity.**)

F is convex, twice continuously differentiable and L -smooth.

Assumption 7 (Smoothness and convexity.)

F is convex, twice continuously differentiable and L -smooth.

Theorem 4 (Convergence of MCM, convex case)

Under A1, A3, A7, for K in \mathbb{N} , with a large enough step-size $\gamma = \sqrt{\frac{\delta_0^2 Nb}{(\omega_{\text{up}}+1)\sigma^2 K}}$, denoting $\bar{w}_K = \frac{1}{K} \sum_{i=0}^{K-1} w_i$, we have:

$$\mathbb{E}[F(\bar{w}_K) - F_*] \leq 2\sqrt{\frac{\delta_0^2 (\omega_{\text{up}} + 1)\sigma^2}{NbK}} + O\left(\frac{\omega_{\text{up}}\omega_{\text{down}}}{K}\right).$$

Assumption 7 (Smoothness and convexity.)

F is convex, twice continuously differentiable and L -smooth.

Theorem 4 (Convergence of MCM, convex case)

Under A1, A3, A7, for K in \mathbb{N} , with a large enough step-size $\gamma = \sqrt{\frac{\delta_0^2 Nb}{(\omega_{\text{up}}+1)\sigma^2 K}}$, denoting $\bar{w}_K = \frac{1}{K} \sum_{i=0}^{K-1} w_i$, we have:

$$\mathbb{E}[F(\bar{w}_K) - F_*] \leq 2 \underbrace{\sqrt{\frac{\delta_0^2 (\omega_{\text{up}} + 1) \sigma^2}{NbK}}}_{\text{dominant term}} + \underbrace{O\left(\frac{\omega_{\text{up}} \omega_{\text{dwn}}}{K}\right)}_{\text{lower order term}}.$$

Assumption 7 (Smoothness and convexity.)

F is convex, twice continuously differentiable and L -smooth.

Theorem 4 (Convergence of MCM, convex case)

Under A1, A3, A7, for K in \mathbb{N} , with a large enough step-size $\gamma = \sqrt{\frac{\delta_0^2 Nb}{(\omega_{\text{up}}+1)\sigma^2 K}}$, denoting $\bar{w}_K = \frac{1}{K} \sum_{i=0}^{K-1} w_i$, we have:

$$\mathbb{E}[F(\bar{w}_K) - F_*] \leq \underbrace{2\sqrt{\frac{\delta_0^2 (\omega_{\text{up}}+1)\sigma^2}{NbK}}}_{\text{dominant term}} + \underbrace{O\left(\frac{\omega_{\text{up}}\omega_{\text{down}}}{K}\right)}_{\text{lower order term}}.$$

- independent of ω_{down}
- identical to Diana (uni-compression)
- depends on ω_{down}
- asymptotically negligible

Assumption 7 (Smoothness and convexity.)

F is convex, twice continuously differentiable and L -smooth.

Theorem 4 (Convergence of MCM, convex case)

Under A1, A3, A7, for K in \mathbb{N} , with a large enough step-size $\gamma = \sqrt{\frac{\delta_0^2 Nb}{(\omega_{\text{up}}+1)\sigma^2 K}}$, denoting $\bar{w}_K = \frac{1}{K} \sum_{i=0}^{K-1} w_i$, we have:

$$\mathbb{E}[F(\bar{w}_K) - F_*] \leq 2\sqrt{\frac{\delta_0^2 (\omega_{\text{up}}+1)\sigma^2}{NbK}} + O\left(\frac{\omega_{\text{up}}\omega_{\text{dwn}}}{K}\right).$$

Moreover if $\sigma^2 = 0$, we recover a faster convergence:

$$\mathbb{E}[F(\bar{w}_K) - F_*] = O(K^{-1}).$$

Assumption 7 (Smoothness and convexity.)

F is convex, twice continuously differentiable and L -smooth.

Theorem 4 (Convergence of MCM, convex case)

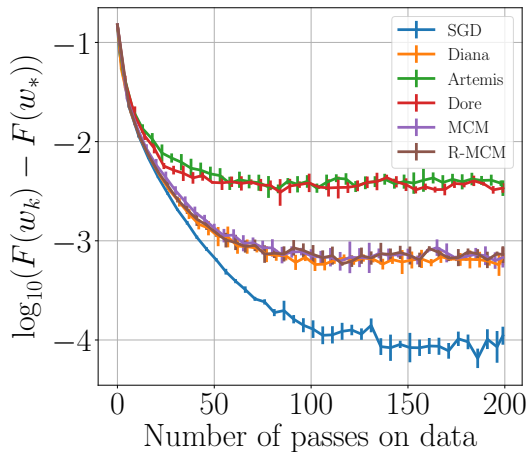
Under A1, A3, A7, for K in \mathbb{N} , with a large enough step-size $\gamma = \sqrt{\frac{\delta_0^2 Nb}{(\omega_{\text{up}}+1)\sigma^2 K}}$, denoting $\bar{w}_K = \frac{1}{K} \sum_{i=0}^{K-1} w_i$, we have:

$$\mathbb{E}[F(\bar{w}_K) - F_*] \leq 2\sqrt{\frac{\delta_0^2 (\omega_{\text{up}} + 1)\sigma^2}{NbK}} + O\left(\frac{\omega_{\text{up}}\omega_{\text{dwn}}}{K}\right).$$

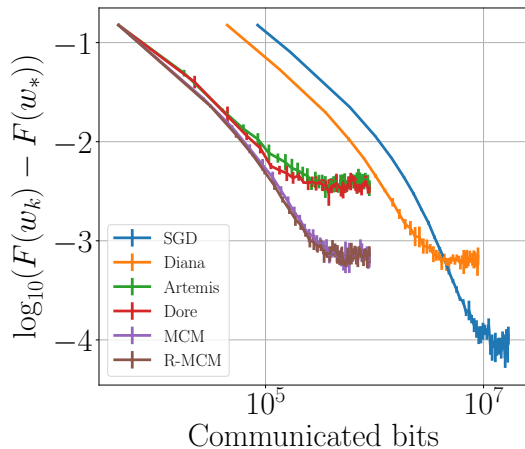
Moreover if $\sigma^2 = 0$, we recover a faster convergence:

$$\mathbb{E}[F(\bar{w}_K) - F_*] = O(K^{-1}).$$

Remark: this result is also extended to both strongly-convex and non-convex cases.



(a) X axis in # iterations



(b) X axis in # bits

Figure 7: Quantum with $b = 400$, $\gamma = 1/L$ (Logistic regression).

Nonconvex framework	MNIST (CNN, $d=2e4$, 4 bits-quantization with norm 2)	Fashion MNIST (FashionSimpleNet, $d=4e5$, 4 bits-quantization with norm 2)	Heterogeneous EMNIST (CNN, $d=2e4$, 4 bits-quantization with norm 2)	CIFAR-10 (LeNet, $d=62e3$, 16 bits-quantization with norm 2)
Accuracy after 300 epochs	SGD: 99.0%	SGD: 92.4%	SGD: 99.0%	SGD: 69.1%
	Diana: 98.9%	Diana: 92.4%	Diana: 98.9%	Diana: 64.0%
	MCM: 98.8%	MCM: 90.6%	MCM: 98.9%	MCM: 63.5%
	Artemis: 97.9%	Artemis: 86.7%	Artemis: 98.3%	Artemis: 54.8%
	Dore: 97.9%	Dore: 87.9%	Dore: 98.5%	Dore: 56.3%
Train loss after 300 epochs	SGD: 0.025	SGD: 0.093	SGD: 0.026	SGD: 0.909
	Diana: 0.034	Diana: 0.141	Diana: 0.031	Diana: 1.047
	MCM: 0.033	MCM: 0.209	MCM: 0.030	MCM: 1.096
	Artemis: 0.075	Artemis: 0.332	Artemis: 0.052	Artemis: 1.342
	Dore: 0.072	Dore: 0.300	Dore: 0.048	Dore: 1.292

Take-away 4

- New algorithm to perform **bidirectional compression**.
- Asymptotically same rate of convergence than **unidirectional compression**.

Take-away 5

- Local gradients computed on a **“perturbed model”** (more challenging).

Additional contributions of the article:

- Randomized-MCM with independent compressions: improves convergence in the quadratic case.

III. Beyond worst-case analysis

↳ To limit the number of bits exchanged, we **compress** the uplink signal before transmitting it.
Big question: what is the impact of C on convergence?

↳ To limit the number of bits exchanged, we **compress** the uplink signal before transmitting it.

Big question: what is the impact of \mathcal{C} on convergence?

Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \frac{\gamma}{N} \sum_{i=1}^N \mathcal{C}(g_k^i(w_{k-1}))$.

Assumption

There exists a constant $\omega \in \mathbb{R}_+^*$ s.t. \mathcal{C} satisfies, for all z in \mathbb{R}^d :

$$\mathbb{E}[\mathcal{C}(z)] = z \quad \text{and} \quad \mathbb{E}[\|\mathcal{C}(z) - z\|^2] \leq \omega \|z\|^2 .$$

↪ To limit the number of bits exchanged, we **compress** the uplink signal before transmitting it.

Big question: what is the impact of \mathcal{C} on convergence?

Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \frac{\gamma}{N} \sum_{i=1}^N \mathcal{C}(g_k^i(w_{k-1}))$.

Assumption

There exists a constant $\omega \in \mathbb{R}_+^*$ s.t. \mathcal{C} satisfies, for all z in \mathbb{R}^d :

$$\mathbb{E}[\mathcal{C}(z)] = z \quad \text{and} \quad \mathbb{E}[\|\mathcal{C}(z) - z\|^2] \leq \omega \|z\|^2.$$

- To go beyond this *worst-case* assumption and provide a tighter analyse.
- Focus on the LSR framework, which is popular for fine-grained analyses.

↪ To limit the number of bits exchanged, we **compress** the uplink signal before transmitting it.

Big question: what is the impact of \mathcal{C} on convergence?

Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \frac{\gamma}{N} \sum_{i=1}^N \mathcal{C}(g_k^i(w_{k-1}))$.

Assumption

There exists a constant $\omega \in \mathbb{R}_+^*$ s.t. \mathcal{C} satisfies, for all z in \mathbb{R}^d :

$$\mathbb{E}[\mathcal{C}(z)] = z \quad \text{and} \quad \mathbb{E}[\|\mathcal{C}(z) - z\|^2] \leq \omega \|z\|^2.$$

- To go beyond this *worst-case* assumption and provide a tighter analyse.
- Focus on the LSR framework, which is popular for fine-grained analyses.

Final goal: highlight the differences in convergence between several unbiased compression schemes having the *same* variance increase.

↳ To limit the number of bits exchanged, we **compress** the uplink signal before transmitting it.

Big question: what is the impact of \mathcal{C} on convergence?

Compressed distributed SGD: $\forall k \in \mathbb{N}, w_k = w_{k-1} - \frac{\gamma}{N} \sum_{i=1}^N \mathcal{C}(g_k^i(w_{k-1}))$.

Assumption

There exists a constant $\omega \in \mathbb{R}_+^*$ s.t. \mathcal{C} satisfies, for all z in \mathbb{R}^d :

$$\mathbb{E}[\mathcal{C}(z)] = z \quad \text{and} \quad \mathbb{E}[\|\mathcal{C}(z) - z\|^2] \leq \omega \|z\|^2.$$

- To go beyond this *worst-case* assumption and provide a tighter analyse.
- Focus on the LSR framework, which is popular for fine-grained analyses.

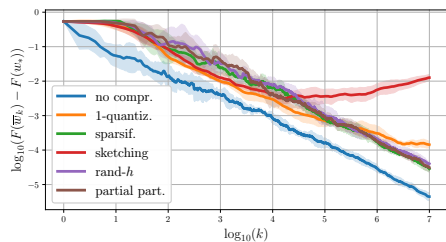
Simplified setting for this presentation:

- $N = 1$ client.
- The client accesses K i.i.d. observations $(x_k, y_k)_{k \in \{1, \dots, K\}} \sim \mathcal{D}^{\otimes K}$, such that there exists a well-defined model w_* :

$$\forall k \in \{1, \dots, K\}, \quad y_k = \langle x_k, w_* \rangle + \varepsilon_k^i, \quad \text{with } \varepsilon_k \sim \mathcal{N}(0, \sigma^2).$$

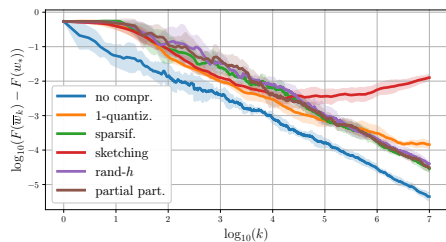
5 compressors: 4 scenarios, 4 different behaviors.

5 compressors: 4 scenarios, 4 different behaviors.

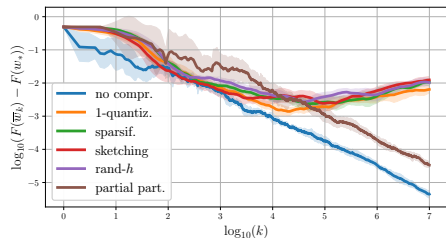


Sketching is very bad, quantiz. is slightly worse.

5 compressors: 4 scenarios, 4 different behaviors.



Sketching is very bad, quantiz. is slightly worse.

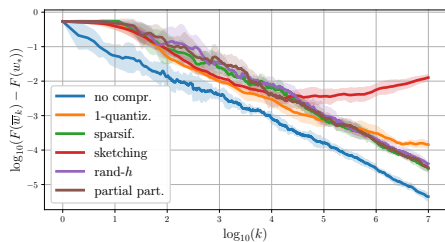


Only partial part. is good.

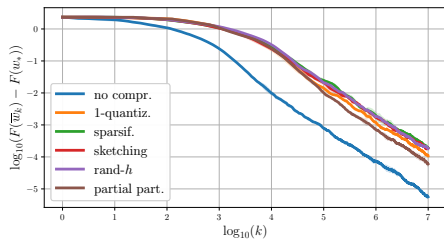
Comparing various compressors in different scenarios



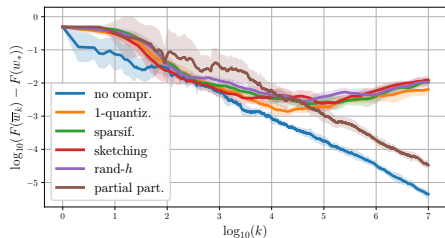
5 compressors: 4 scenarios, 4 different behaviors.



Sketching is very bad, quantiz. is slightly worse.



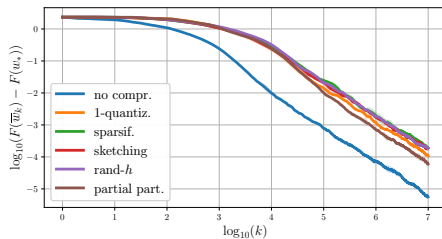
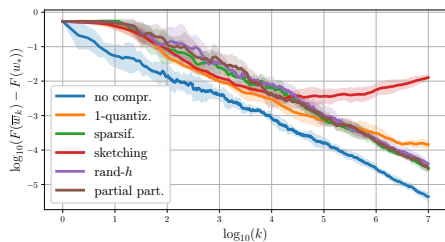
All compressors are equivalent and behave well.



Only partial part. is good.

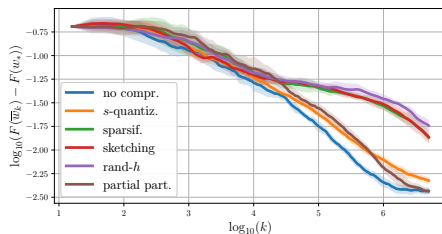
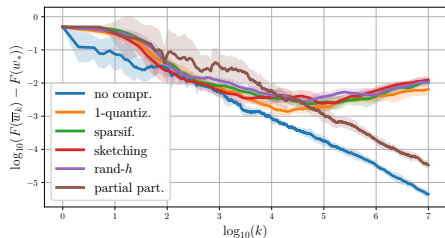
Comparing various compressors in different scenarios

5 compressors: 4 scenarios, 4 different behaviors.



Sketching is very bad, quantiz. is slightly worse.

All compressors are equivalent and behave well.

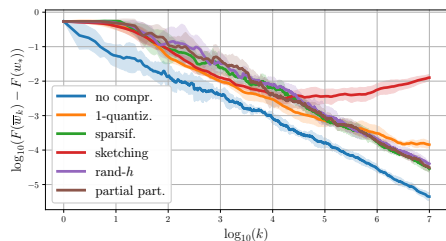


Only partial part. is good.

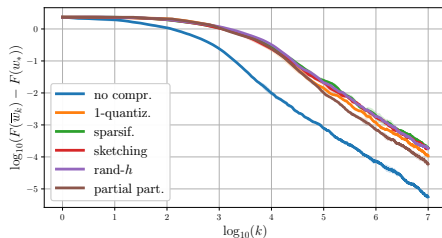
Quantiz. and partial part. are good.

Comparing various compressors in different scenarios

5 compressors: 4 scenarios, 4 different behaviors.

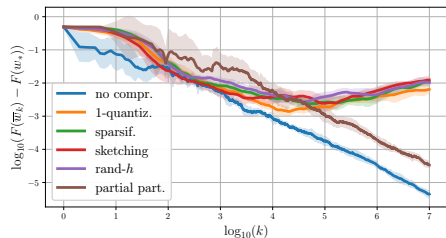


Can we explain this four different behaviors?

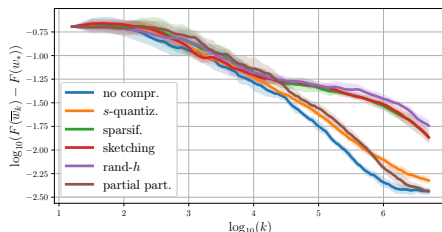


Sketching is very bad, quantiz. is slightly worse.

All compressors are equivalent and behave well.



Only partial part. is good.



Quantiz. and partial part. are good.

Definition 1 (Linear Stochastic Approximation, LSA)

Let $w_0 \in \mathbb{R}^d$ be the initialization, the linear stochastic approximation¹ recursion is defined as:

$$w_k = w_{k-1} - \gamma \nabla F(w_{k-1}) + \gamma \xi_k (w_{k-1} - w_*), \quad k \in \mathbb{N}, \quad (\text{LSA})$$

- $\gamma > 0$: step size,
- $(\xi_k)_{k \in \mathbb{N}^*}$: sequence of i.i.d. zero-centered random fields that characterizes the stochastic oracle on $\nabla F(\cdot)$.

¹While in LSA literature, both the mean-field ∇F and the noise-field (ξ_k) are linear, we do not here consider the noise fields to be linear.

Definition 1 (Linear Stochastic Approximation, LSA)

Let $w_0 \in \mathbb{R}^d$ be the initialization, the linear stochastic approximation¹ recursion is defined as:

$$w_k = w_{k-1} - \gamma \nabla F(w_{k-1}) + \gamma \xi_k (w_{k-1} - w_*), \quad k \in \mathbb{N}, \quad (\text{LSA})$$

- $\gamma > 0$: step size,
- $(\xi_k)_{k \in \mathbb{N}^*}$: sequence of i.i.d. zero-centered random fields that characterizes the stochastic oracle on $\nabla F(\cdot)$.

We assume F quadratic:

- H_F : its Hessian
- μ : its smallest eigenvalue.

For any k in \mathbb{N} , with $\eta_k = w_k - w_*$, we get equivalently:

$$\eta_k = (I - \gamma H_F) \eta_{k-1} + \gamma \xi_k (\eta_{k-1}), \quad k \in \mathbb{N}.$$

¹While in LSA literature, both the mean-field ∇F and the noise-field (ξ_k) are linear, we do not here consider the noise fields to be linear.

Algorithm 1 (LMS with a single worker)

We have for all $k \in \mathbb{N}$:

$$w_k = w_{k-1} - \gamma(\langle w_{k-1}, x_k \rangle - y_k)x_k,$$

Equivalently, for $w \in \mathbb{R}^d$:

$$\xi_k(w) = (x_k x_k^\top - \mathbb{E}[x_1 x_1^\top])w + (\langle w_*, x_k \rangle - y_k)x_k.$$

Algorithm 1 (LMS with a single worker)

We have for all $k \in \mathbb{N}$:

$$w_k = w_{k-1} - \gamma(\langle w_{k-1}, x_k \rangle - y_k)x_k,$$

Equivalently, for $w \in \mathbb{R}^d$:

$$\xi_k(w) = (x_k x_k^\top - \mathbb{E}[x_1 x_1^\top])w + (\langle w_*, x_k \rangle - y_k)x_k.$$

Algorithm 2 (Centralized compressed LMS)

At any step k in $\{1, \dots, K\}$, we have an oracle $g_k(\cdot)$ of the gradient of the objective function F and a random compression mechanism $\mathcal{C}_k(\cdot)$.

For any step-size $\gamma > 0$ and any $k \in \mathbb{N}^*$, the resulting sequence of iterates $(w_k)_{k \in \mathbb{N}}$ satisfies:

$$w_k = w_{k-1} - \gamma \mathcal{C}_k(g_k(w_{k-1})).$$

Equivalently, for $w \in \mathbb{R}^d$:

$$\xi_k(w) = \nabla F(w) - \mathcal{C}_k(g_k(w)).$$

Algorithm 1 (LMS with a single worker)

We have for all $k \in \mathbb{N}$:

$$w_k = w_{k-1} - \gamma(\langle w_{k-1}, x_k \rangle - y_k)x_k,$$

Equivalently, for $w \in \mathbb{R}^d$:

$$\xi_k(w) = (x_k x_k^\top - \mathbb{E}[x_1 x_1^\top])w + (\langle w_*, x_k \rangle - y_k)x_k.$$

Algorithm 2 (Centralized compressed LMS)

At any step k in $\{1, \dots, K\}$, we have an oracle $g_k(\cdot)$ of the gradient of the objective function F and a random compression mechanism $\mathcal{C}_k(\cdot)$.

For any step-size $\gamma > 0$ and any $k \in \mathbb{N}^*$, the resulting sequence of iterates $(w_k)_{k \in \mathbb{N}}$ satisfies:

$$w_k = w_{k-1} - \gamma \mathcal{C}_k(g_k(w_{k-1})).$$

Equivalently, for $w \in \mathbb{R}^d$:

$$\xi_k(w) = \nabla F(w) - \mathcal{C}_k(g_k(w)).$$

Most analyses of (LSA)

[Blu54, Lju77, LS83] assume either:

1. The field ξ_k is either linear [see KT03, BMP12, LP21] i.e. for any $z, z' \in \mathbb{R}^d$,

$$\xi_k(z) - \xi_k(z') = \xi_k(z - z').$$

2. The noise-field is Lipschitz in squared expectation [MB11, Bac14, DDB20, GP23].
i.e. for any $z, z' \in \mathbb{R}^d$

$$\mathbb{E}[\|\xi_k(z) - \xi_k(z')\|^2] \leq C\|z - z'\|^2.$$

Algorithm 1 (LMS with a single worker)

We have for all $k \in \mathbb{N}$:

$$w_k = w_{k-1} - \gamma(\langle w_{k-1}, x_k \rangle - y_k)x_k,$$

Equivalently, for $w \in \mathbb{R}^d$:

$$\xi_k(w) = (x_k x_k^\top - \mathbb{E}[x_1 x_1^\top])w + (\langle w_*, x_k \rangle - y_k)x_k.$$

Algorithm 2 (Centralized compressed LMS)

At any step k in $\{1, \dots, K\}$, we have an oracle $g_k(\cdot)$ of the gradient of the objective function F and a random compression mechanism $\mathcal{C}_k(\cdot)$.

For any step-size $\gamma > 0$ and any $k \in \mathbb{N}^*$, the resulting sequence of iterates $(w_k)_{k \in \mathbb{N}}$ satisfies:

$$w_k = w_{k-1} - \gamma \mathcal{C}_k(g_k(w_{k-1})).$$

Equivalently, for $w \in \mathbb{R}^d$:

$$\xi_k(w) = \nabla F(w) - \mathcal{C}_k(g_k(w)).$$

Most analyses of (LSA)

[Blu54, Lju77, LS83] assume either:

1. The field ξ_k is either linear [see KT03, BMP12, LP21] i.e. for any $z, z' \in \mathbb{R}^d$,

$$\xi_k(z) - \xi_k(z') = \xi_k(z - z').$$

2. The noise-field is Lipschitz in squared expectation [MB11, Bac14, DDB20, GP23].
i.e. for any $z, z' \in \mathbb{R}^d$

$$\mathbb{E}[\|\xi_k(z) - \xi_k(z')\|^2] \leq C\|z - z'\|^2.$$

\implies Specificity and bottleneck of compression: the resulting field **does not** satisfy such assumptions.

Definition 2 (Additive and multiplicative noise)

Under the setting of (LSA), for any k in \mathbb{N}^ :*

$$\xi_k^{\text{add}} := \xi_k(0) \quad \text{and} \quad \xi_k^{\text{mult}} : z \in \mathbb{R}^d \mapsto \xi_k(z) - \xi_k^{\text{add}}.$$

Definition 2 (Additive and multiplicative noise)

Under the setting of (LSA), for any k in \mathbb{N}^* :

$$\xi_k^{\text{add}} := \xi_k(0) \quad \text{and} \quad \xi_k^{\text{mult}} : z \in \mathbb{R}^d \mapsto \xi_k(z) - \xi_k^{\text{add}}.$$

Assumption (Second moment of the multiplicative noise)

$\exists \mathcal{M}_1, \mathcal{M}_2 > 0$ s.t. for any η in \mathbb{R}^d :

1. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq 2\mathcal{M}_2 \|H_F^{1/2} \eta\|^2 + 4\mathcal{A}.$
2. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq \mathcal{M}_1 \|H_F^{1/2} \eta\| + 3\mathcal{M}_2 \|H_F^{1/2} \eta\|^2.$

Definition 2 (Additive and multiplicative noise)

Under the setting of (LSA), for any k in \mathbb{N}^* :

$$\xi_k^{\text{add}} := \xi_k(0) \quad \text{and} \quad \xi_k^{\text{mult}} : z \in \mathbb{R}^d \mapsto \xi_k(z) - \xi_k^{\text{add}}.$$

Assumption (Second moment of the multiplicative noise)

$\exists \mathcal{M}_1, \mathcal{M}_2 > 0$ s.t. for any η in \mathbb{R}^d :

1. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq 2\mathcal{M}_2 \|H_F^{1/2} \eta\|^2 + 4\mathcal{A}.$
2. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq \mathcal{M}_1 \|H_F^{1/2} \eta\| + 3\mathcal{M}_2 \|H_F^{1/2} \eta\|^2.$

Classical assumption

Hölder-type assumption
(new)

Definition 2 (Additive and multiplicative noise)

Under the setting of (LSA), for any k in \mathbb{N}^* :

$$\xi_k^{\text{add}} := \xi_k(0) \quad \text{and} \quad \xi_k^{\text{mult}} : z \in \mathbb{R}^d \mapsto \xi_k(z) - \xi_k^{\text{add}}.$$

Assumption (Second moment of the multiplicative noise)

$\exists \mathcal{M}_1, \mathcal{M}_2 > 0$ s.t. for any η in \mathbb{R}^d :

1. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq 2\mathcal{M}_2 \|H_F^{1/2} \eta\|^2 + 4\mathcal{A}.$
2. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq \mathcal{M}_1 \|H_F^{1/2} \eta\| + 3\mathcal{M}_2 \|H_F^{1/2} \eta\|^2.$

Classical assumption

Hölder-type assumption (new)

$\mathcal{M}_1 = 0$ if the random field is linear,
 $\mathcal{M}_1 \neq 0$ for quantization

Definition 2 (Additive and multiplicative noise)

Under the setting of (LSA), for any k in \mathbb{N}^* :

$$\xi_k^{\text{add}} := \xi_k(0) \quad \text{and} \quad \xi_k^{\text{mult}} : z \in \mathbb{R}^d \mapsto \xi_k(z) - \xi_k^{\text{add}}.$$

Assumption (Second moment of the multiplicative noise)

$\exists \mathcal{M}_1, \mathcal{M}_2 > 0$ s.t. for any η in \mathbb{R}^d :

1. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq 2\mathcal{M}_2 \|H_F^{1/2} \eta\|^2 + 4\mathcal{A}.$
2. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq \mathcal{M}_1 \|H_F^{1/2} \eta\| + 3\mathcal{M}_2 \|H_F^{1/2} \eta\|^2.$

Classical assumption

Hölder-type assumption (new)

$\mathcal{M}_1 = 0$ if the random field is linear,
 $\mathcal{M}_1 \neq 0$ for quantization because:

$$\mathbb{E}[\|\mathcal{C}(z) - \mathcal{C}(z')\|^2] \leq 12\sqrt{d} \min(\|z\|, \|z'\|) \|z - z'\| + 3(\omega + 1) \|z - z'\|^2$$

Definition 2 (Additive and multiplicative noise)

Under the setting of (LSA), for any k in \mathbb{N}^* :

$$\xi_k^{\text{add}} := \xi_k(0) \quad \text{and} \quad \xi_k^{\text{mult}} : z \in \mathbb{R}^d \mapsto \xi_k(z) - \xi_k^{\text{add}}.$$

Assumption (Second moment of the multiplicative noise)

$\exists \mathcal{M}_1, \mathcal{M}_2 > 0$ s.t. for any η in \mathbb{R}^d :

1. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq 2\mathcal{M}_2 \|H_F^{1/2} \eta\|^2 + 4\mathcal{A}.$
2. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq \mathcal{M}_1 \|H_F^{1/2} \eta\| + 3\mathcal{M}_2 \|H_F^{1/2} \eta\|^2.$

Definition 3 (Ania's covariance.)

Under (LSA), we define the covariance of the additive noise: $\mathfrak{C}_{\text{ania}} = \mathbb{E}[\xi_1^{\text{add}} \otimes \xi_1^{\text{add}}].$

Definition 2 (Additive and multiplicative noise)

Under the setting of (LSA), for any k in \mathbb{N}^* :

$$\xi_k^{\text{add}} := \xi_k(0) \quad \text{and} \quad \xi_k^{\text{mult}} : z \in \mathbb{R}^d \mapsto \xi_k(z) - \xi_k^{\text{add}}.$$

Assumption (Second moment of the multiplicative noise)

$\exists \mathcal{M}_1, \mathcal{M}_2 > 0$ s.t. for any η in \mathbb{R}^d :

1. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq 2\mathcal{M}_2 \|H_F^{1/2} \eta\|^2 + 4\mathcal{A}.$
2. $\mathbb{E}[\|\xi_1^{\text{mult}}(\eta)\|^2] \leq \mathcal{M}_1 \|H_F^{1/2} \eta\| + 3\mathcal{M}_2 \|H_F^{1/2} \eta\|^2.$

Definition 3 (Ania's covariance.)

Under (LSA), we define the covariance of the additive noise: $\mathfrak{C}_{\text{ania}} = \mathbb{E}[\xi_1^{\text{add}} \otimes \xi_1^{\text{add}}].$

Theorem 5 (Asymptotic result, from [PJ92])

Under some assumptions. Consider a sequence $(w_k)_{k \in \mathbb{N}^*}$ produced in the setting of (LSA) for a step-size $(\gamma_K)_{K \in \mathbb{N}^*}$ s.t. $\gamma_K = 1/\sqrt{K}$. Then we have:

$$\sqrt{K}(\bar{w}_K - w_*) \xrightarrow[K \rightarrow +\infty]{\mathcal{L}} \mathcal{N}(0, H_F^{-1} \mathfrak{C}_{\text{ania}} H_F^{-1}).$$

Theorem 6 (“Non-asymptotic convergence rate”)

Under some assumptions. Consider a sequence $(w_k)_{k \in \mathbb{N}^}$ produced by the setting of (LSA), for a constant step-size γ verifying some assumptions. Then for any horizon K , we have*

$$\mathbb{E}[F(\bar{w}_{K-1}) - F(w_*)] \leq \frac{1}{2K} \left(\min \left(\frac{\|H_F^{-1/2} \eta_0\|}{\gamma \sqrt{K}}, \frac{\|\eta_0\|}{\sqrt{\gamma}} \right) + \sqrt{\text{Tr}(\mathfrak{C}_{\text{ania}} H_F^{-1})} + O(\mu^{-1/2} \gamma^{1/4}) \right)^2.$$

Theorem 6 (“Non-asymptotic convergence rate”)

Under some assumptions. Consider a sequence $(w_k)_{k \in \mathbb{N}^*}$ produced by the setting of (LSA), for a constant step-size γ verifying some assumptions. Then for any horizon K , we have

$$\mathbb{E}[F(\bar{w}_{K-1}) - F(w_*)] \leq \frac{1}{2K} \left(\min \left(\frac{\|H_F^{-1/2} \eta_0\|}{\gamma \sqrt{K}}, \frac{\|\eta_0\|}{\sqrt{\gamma}} \right) + \sqrt{\text{Tr}(\mathfrak{C}_{\text{ania}} H_F^{-1}) + O(\mu^{-1/2} \gamma^{1/4})} \right)^2.$$

Bias term, as in [BM13, DB15]

classical asymptotic noise term in CLT for (LSA)

asymptotically negligible for $\gamma = o(1)$, comes from multiplicative noise

$$\eta_k = w_k - w_*$$

$\mathfrak{C}_{\text{ania}}$: additive noise's covariance

H_F : Hessian

$$\mu = \min(\text{eig}(H_F))$$

Theorem 6 (“Non-asymptotic convergence rate”)

Under some assumptions. Consider a sequence $(w_k)_{k \in \mathbb{N}^*}$ produced by the setting of (LSA), for a constant step-size γ verifying some assumptions. Then for any horizon K , we have

$$\mathbb{E}[F(\bar{w}_{K-1}) - F(w_*)] \leq \frac{1}{2K} \left(\min \left(\frac{\|H_F^{-1/2} \eta_0\|}{\gamma \sqrt{K}}, \frac{\|\eta_0\|}{\sqrt{\gamma}} \right) + \sqrt{\text{Tr}(\mathfrak{C}_{\text{ania}} H_F^{-1}) + O(\mu^{-1/2} \gamma^{1/4})} \right)^2.$$

Bias term, as in [BM13, DB15]

classical asymptotic noise term in CLT for (LSA)

asymptotically negligible for $\gamma = o(1)$,
comes from multiplicative noise

Remarks:

- Asymptotically, the dominant term is $\sqrt{\text{Tr}(\mathfrak{C}_{\text{ania}} H_F^{-1})}$.
- Contrary to [BM13], the convergence rate *is not* necessarily independent of μ .
- Examining the explicit formulas of $\text{Tr}(\mathfrak{C}_{\text{ania}} H_F^{-1})$ allows to determine the convergence rate.

$$\eta_k = w_k - w_*$$

$\mathfrak{C}_{\text{ania}}$: additive noise's covariance

H_F : Hessian

$$\mu = \min(\text{eig}(H_F))$$

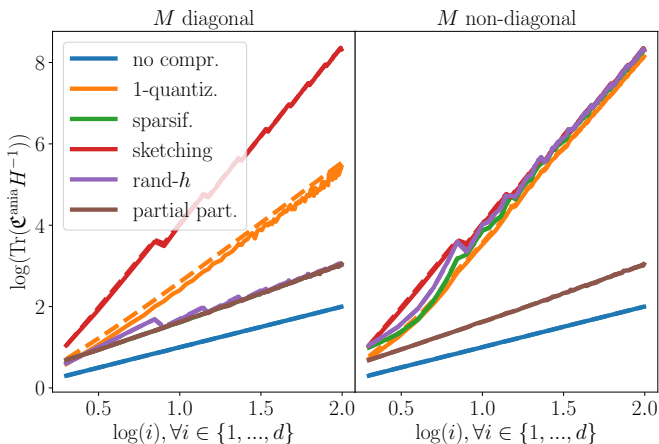
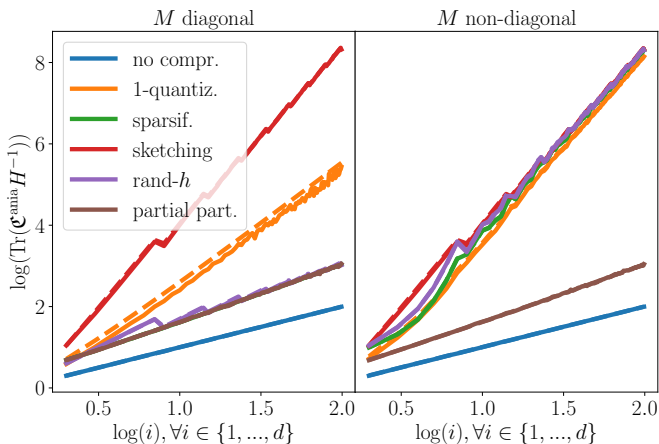


Figure 8: $\text{Tr}(\mathbf{C}_{\text{ania}} H^{-1})$ - $K = 10^3, d \in \llbracket 2, 100 \rrbracket, D = \text{Diag}((1/i^4)_{i=1}^d)$. Left: H diagonal. Right: H non-diagonal. (Plain line: empirical values; dashed lines: theoretical)

$\forall k \in \{1, \dots, K\}, x_k \sim \mathcal{N}(0, H)$, with $H = QDQ^T$, $D = \text{Diag}((1/i^4)_{i=1}^d)$ and Q an orthogonal matrix.



Depending on the compression scheme:

Classical LMS: $\mathcal{C}_{\text{ania}} = H \quad (\times \sigma^2)$

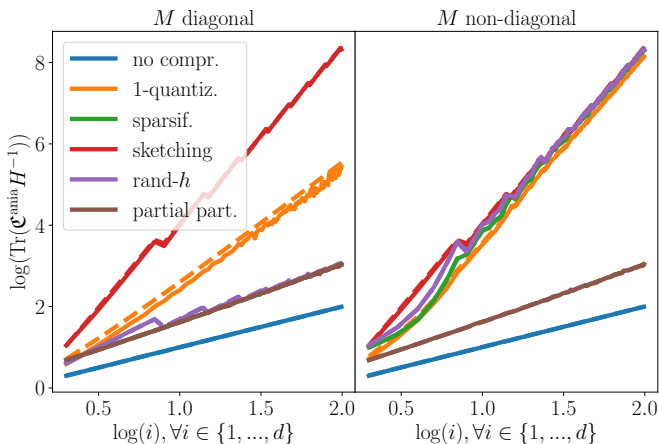
Partial part.: $\mathcal{C}_{\text{ania}} = aH$

Sparsification: $\mathcal{C}_{\text{ania}} = a'H + b\text{Diag}(H)$

Sketching: $\mathcal{C}_{\text{ania}} = a''H + b'\text{Tr}(H)I_d$

Figure 8: $\text{Tr}(\mathcal{C}_{\text{ania}}H^{-1})$ - $K = 10^3, d \in \llbracket 2, 100 \rrbracket, D = \text{Diag}((1/i^4)_{i=1}^d)$. Left: H diagonal. Right: H non-diagonal. (Plain line: empirical values; dashed lines: theoretical)

$\forall k \in \{1, \dots, K\}, x_k \sim \mathcal{N}(0, H)$, with $H = QDQ^T$, $D = \text{Diag}((1/i^4)_{i=1}^d)$ and Q an orthogonal matrix.



Depending on the compression scheme:

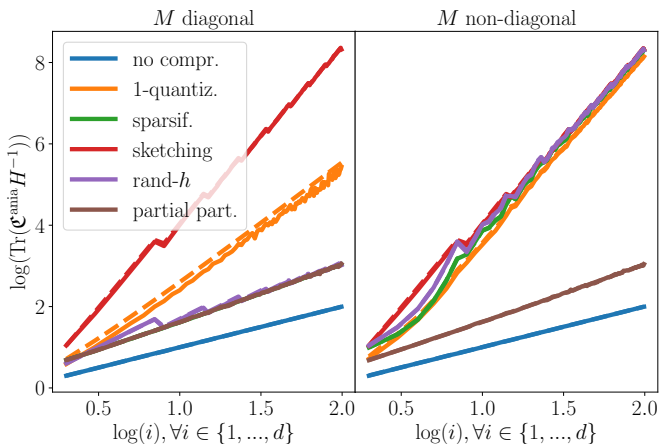
- Classical LMS: $\mathcal{C}_{\text{ania}} = H \quad (\times \sigma^2)$
- Partial part.: $\mathcal{C}_{\text{ania}} = aH$
- Sparsification: $\mathcal{C}_{\text{ania}} = a'H + b\text{Diag}(H)$
- Sketching: $\mathcal{C}_{\text{ania}} = a''H + b'\text{Tr}(H)I_d$

Structured noise

Isotropic noise

Figure 8: $\text{Tr}(\mathcal{C}_{\text{ania}} H^{-1})$ - $K = 10^3, d \in \llbracket 2, 100 \rrbracket, D = \text{Diag}((1/i^4)_{i=1}^d)$. Left: H diagonal. Right: H non-diagonal. (Plain line: empirical values; dashed lines: theoretical)

$\forall k \in \{1, \dots, K\}, x_k \sim \mathcal{N}(0, H)$, with $H = QDQ^T$, $D = \text{Diag}((1/i^4)_{i=1}^d)$ and Q an orthogonal matrix.



Depending on the compression scheme:

Classical LMS: $\mathfrak{C}_{\text{ania}} = H \quad (\times \sigma^2)$

Partial part.: $\mathfrak{C}_{\text{ania}} = aH$

Sparsification: $\mathfrak{C}_{\text{ania}} = a'H + b\text{Diag}(H)$

Sketching: $\mathfrak{C}_{\text{ania}} = a''H + b'\text{Tr}(H)I_d$

Structured noise

Isotropic noise

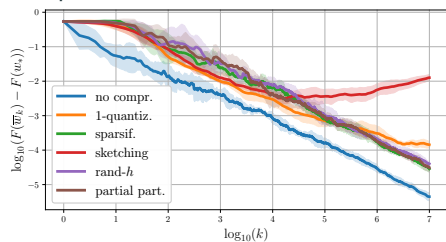
- Significantly impacts the limit distribution with a rate proportional to $\text{Tr}(H^{-1})$.
- Same variance but different behaviors!

Figure 8: $\text{Tr}(\mathfrak{C}_{\text{ania}} H^{-1})$ - $K = 10^3, d \in \llbracket 2, 100 \rrbracket, D = \text{Diag}((1/i^4)_{i=1}^d)$. Left: H diagonal. Right: H non-diagonal. (Plain line: empirical values; dashed lines: theoretical)

$\forall k \in \{1, \dots, K\}, x_k \sim \mathcal{N}(0, H)$, with $H = QDQ^T, D = \text{Diag}((1/i^4)_{i=1}^d)$ and Q an orthogonal matrix.

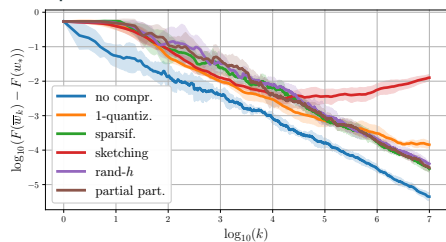
5 compressors: 4 scenarios, 4 different behaviors.

5 compressors: 4 scenarios, 4 different behaviors.

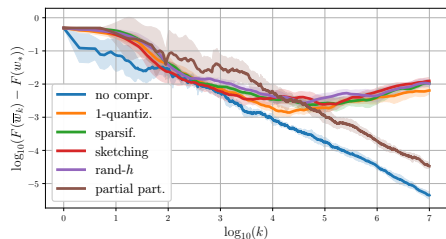


Fast eigenvalues' decay, diagonal covariance H .

5 compressors: 4 scenarios, 4 different behaviors.

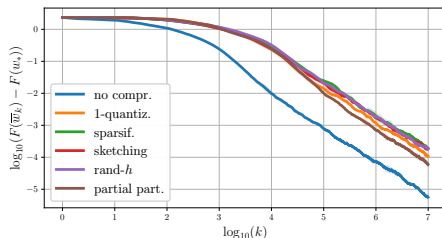
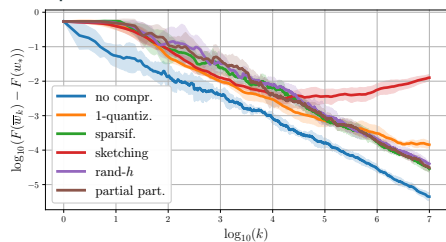


Fast eigenvalues' decay, diagonal covariance H .



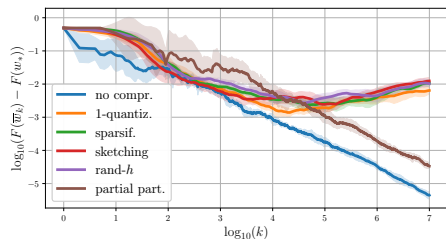
Fast eigenvalues' decay, non-diagonal covariance H .

5 compressors: 4 scenarios, 4 different behaviors.



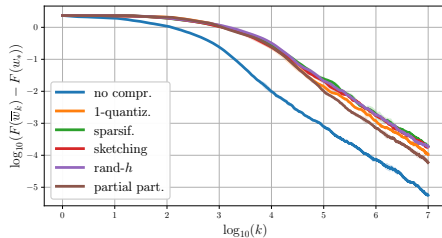
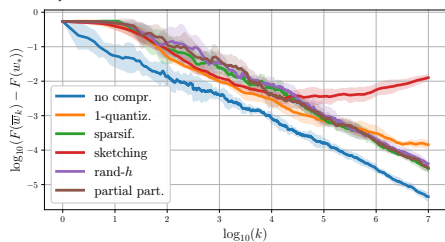
Fast eigenvalues' decay, diagonal covariance H .

Slow eigenvalues' decay, non-diagonal covariance H .

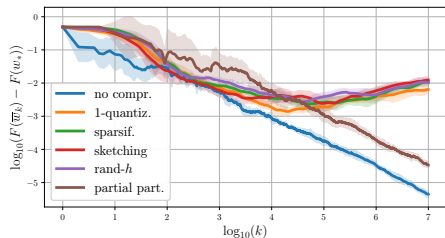


Fast eigenvalues' decay, non-diagonal covariance H .

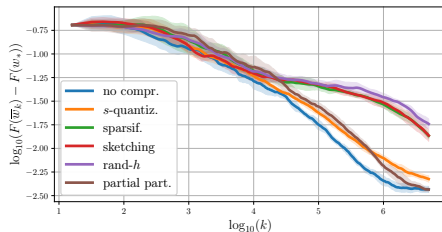
5 compressors: 4 scenarios, 4 different behaviors.



Fast eigenvalues' decay, diagonal covariance H .



Slow eigenvalues' decay, non-diagonal covariance H .



Fast eigenvalues' decay, non-diagonal covariance H .

Cifar10 with standardization (constant diagonal covariance H).

Summary of the contributions of the article:

- Analyze (LSA) under weak regularity assumptions of the noise field $(\xi_k)_k$.
- Provide a non-asymptotic theorem.
- Underline the key impact on convergence of the ania's covariance $\mathcal{C}_{\text{ania}}$.
- Describe the link between, the compressor \mathcal{C} , the features' covariance H and the ania's covariance $\mathcal{C}_{\text{ania}}$.
- Show how to compute the ania's covariance $\mathcal{C}_{\text{ania}}$.
- Study the FL setting with heterogeneous clients.

Summary of the contributions of the article:

- Analyze (LSA) under weak regularity assumptions of the noise field $(\xi_k)_k$.
- Provide a non-asymptotic theorem.
- Underline the key impact on convergence of the ania's covariance $\mathcal{C}_{\text{ania}}$.
- Describe the link between, the compressor \mathcal{C} , the features' covariance H and the ania's covariance $\mathcal{C}_{\text{ania}}$.
- Show how to compute the ania's covariance $\mathcal{C}_{\text{ania}}$.
- Study the FL setting with heterogeneous clients.

Examples of take-aways:

Take-away 6

- *Quantization not Lipschitz in squared expectation but satisfy a **Hölder-type** condition.*
- *Convergence degraded, yet achieve a **rate comparable to projection based compressors**.*

Take-away 7

- *Rand-1 and Partial Participation with probability $(1/d)$: **same variance condition**.*
- *But **PP is more robust** to ill conditioned problem.*

Conclusion

Table 2: Summary of contributions.

	Bi-compr.	Heterogeneity	LSR	
I.	✓	✓		Interaction between compression and heterogeneity
II.	✓		(✓)	Asympt. cancels impact of down compression
III.		(✓)	✓	Beyond worst-case analysis

Table 2: Summary of contributions.

	Bi-compr.	Heterogeneity	LSR	
I.	✓	✓		Interaction between compression and heterogeneity
II.	✓		(✓)	Asympt. cancels impact of down compression
III.		(✓)	✓	Beyond worst-case analysis

- I. **Artemis** **Bidirectional compression** to reduce communication cost.
Key impact of **memory** on the convergence on **non-i.i.d.** data.

Table 2: Summary of contributions.

	Bi-compr.	Heterogeneity	LSR	
I.	✓	✓		Interaction between compression and heterogeneity
II.	✓		(✓)	Asympt. cancels impact of down compression
III.		(✓)	✓	Beyond worst-case analysis

I. **Artemis** **Bidirectional compression** to reduce communication cost.
Key impact of **memory** on the convergence on **non-i.i.d.** data.

II. **MCM** Asympt, same rate of convergence **as unidirectional compression**.
Underlines the importance to **not degrade the global model**.

Table 2: Summary of contributions.

	Bi-compr.	Heterogeneity	LSR	
I.	✓	✓		Interaction between compression and heterogeneity
II.	✓		(✓)	Asympt. cancels impact of down compression
III.		(✓)	✓	Beyond worst-case analysis

- I. **Artemis** **Bidirectional compression** to reduce communication cost.
Key impact of **memory** on the convergence on **non-i.i.d.** data.
- II. **MCM** Asympt, same rate of convergence **as unidirectional compression**.
Underlines the importance to **not degrade the global model**.
- III. Beyond the worst-case analysis of compression.
Analyze of the **compressors' covariance**.
Differences between compressors that have the same variance.

Thank you for your attention.

References

- [AJL⁺23] A. Affouard, A. Joly, J. Lombardo, J. Champ, H. Goeau, M. Chouet, H. Gresse, C. Botella, and P. Bonnet. PI@ntnet automatically identified occurrences. v1.8. PI@ntNet, <https://plantnet.org/>, 2023.
- [Bac14] Francis Bach. Adaptivity of averaged stochastic gradient descent to local strong convexity for logistic regression. *The Journal of Machine Learning Research*, 15(1):595–627, 2014.
- [Blu54] Julius R Blum. Multidimensional stochastic approximation methods. *The Annals of Mathematical Statistics*, pages 737–744, 1954.
- [BM13] Francis Bach and Eric Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate $O(1/n)$. In *Neural Information Processing Systems (NIPS)*, pages –, United States, December 2013.
- [BMP12] Albert Benveniste, Michel Métivier, and Pierre Priouret. *Adaptive algorithms and stochastic approximations*, volume 22. Springer Science & Business Media, 2012.
- [CG23] Chat-GPT3. Conversation with chat-gpt3. <https://chat.openai.com>, 2023. Accessed: 2023-08-16.

- [DB15] Alexandre Defossez and Francis Bach. Averaged Least-Mean-Squares: Bias-Variance Trade-offs and Optimal Sampling Distributions. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, pages 205–213. PMLR, February 2015. ISSN: 1938-7228.
- [DBW12] John C. Duchi, Peter L. Bartlett, and Martin J. Wainwright. Randomized Smoothing for Stochastic Optimization. *SIAM Journal on Optimization*, 22(2):674–701, January 2012. Publisher: Society for Industrial and Applied Mathematics.
- [DDB20] Aymeric Dieuleveut, Alain Durmus, and Francis Bach. Bridging the gap between constant step size stochastic gradient descent and markov chains. *Ann. Statist.*, 48(3):1348–1382, 06 2020.
- [DE23] Dall-E. An old castle on a cloud in a miyazaki style. <https://labs.openai.com/>, 2023. Accessed: 2023-04-01.
- [DFB17] Aymeric Dieuleveut, Nicolas Flammarion, and Francis Bach. Harder, better, faster, stronger convergence rates for least-squares regression. *The Journal of Machine Learning Research*, 18(1):3520–3570, 2017. Publisher: JMLR. org.
- [GLQ⁺19] Robert Mansel Gower, Nicolas Loizou, Xun Qian, Alibek Sailanbayev, Egor Shulgin, and Peter Richtárik. SGD: General Analysis and Improved Rates. In *International Conference on Machine Learning*, pages 5200–5209. PMLR, May 2019. ISSN: 2640-3498.

- [GP23] Sébastien Gadat and Fabien Panloup. Optimal non-asymptotic analysis of the Ruppert–Polyak averaging stochastic algorithm. *Stochastic Processes and their Applications*, 156:312–348, February 2023.
- [KT03] Vijay R Konda and John N Tsitsiklis. Linear stochastic approximation driven by slowly varying markov chains. *Systems & control letters*, 50(2):95–102, 2003.
- [Lju77] Lennart Ljung. Analysis of recursive stochastic algorithms. *IEEE transactions on automatic control*, 22(4):551–575, 1977.
- [LLTY20] Xiaorui Liu, Yao Li, Jiliang Tang, and Ming Yan. A Double Residual Compression Algorithm for Efficient Distributed Learning. In *International Conference on Artificial Intelligence and Statistics*, pages 133–143, June 2020. ISSN: 1938-7228 Section: Machine Learning.
- [LP21] Rémi Leluc and François Portier. SGD with Coordinate Sampling: Theory and Practice. *arXiv:2105.11818 [cs, stat]*, May 2021. arXiv: 2105.11818.
- [LS83] Lennart Ljung and Torsten Söderström. *Theory and practice of recursive identification*. MIT press, 1983.
- [MB11] Eric Moulines and Francis Bach. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. *Advances in neural information processing systems*, 24, 2011.

- [MGTR19] Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed Learning with Compressed Gradient Differences. *arXiv:1901.09269 [cs, math, stat]*, June 2019. arXiv: 1901.09269.
- [MMR⁺17] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Artificial Intelligence and Statistics*, pages 1273–1282. PMLR, April 2017. ISSN: 2640-3498.
- [PD20] Constantin Philippenko and Aymeric Dieuleveut. Artemis: tight convergence guarantees for bidirectional compression in Federated Learning. *arXiv:2006.14591 [cs, stat]*, November 2020. arXiv: 2006.14591.
- [PD21] Constantin Philippenko and Aymeric Dieuleveut. Preserved central model for faster bidirectional compression in distributed settings. *Advances in Neural Information Processing Systems*, 34, 2021.
- [PJ92] Boris Polyak and Anatoli Juditsky. Acceleration of Stochastic Approximation by Averaging. *SIAM Journal on Control and Optimization*, 30:838–855, July 1992.
- [SBB⁺18] Kevin Scaman, Francis Bach, Sebastien Bubeck, Laurent Massoulié, and Yin Tat Lee. Optimal Algorithms for Non-Smooth Distributed Optimization in Networks. *Advances in Neural Information Processing Systems*, 31:2740–2749, 2018.

- [TYL⁺19] Hanlin Tang, Chen Yu, Xiangru Lian, Tong Zhang, and Ji 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.
- [Vil09] C. Villani. *Optimal transport : old and new*. Grundlehren der mathematischen wissenschaften. Springer, Berlin, 2009.
- [ZHK19] Shuai Zheng, Ziyue Huang, and James Kwok. Communication-Efficient Distributed Blockwise Momentum SGD with Error-Feedback. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

- Evaluating the type and degree of heterogeneity within a network of clients.
- Compression and neural network: impact in a non-convex setting.
- New schemas of compression with independant coordinate compression.

Back-up on Artemis

Building non-i.i.d. *and* unbalanced datasets using a TSNE representation.

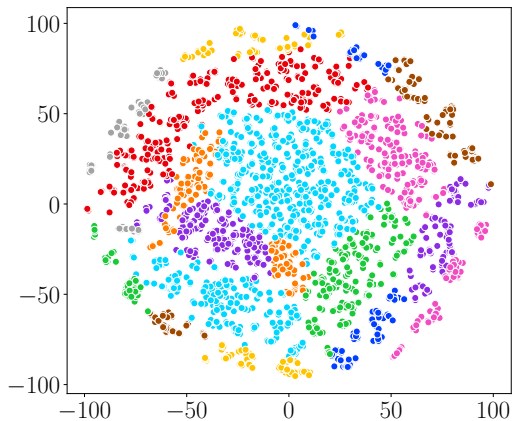


Figure 9: Superconduct

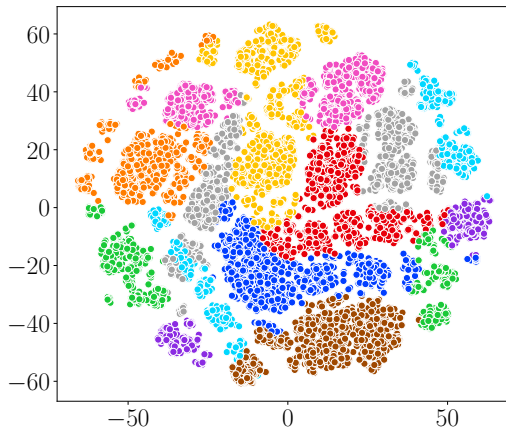


Figure 10: Quantum

We note $\tilde{\mathbf{g}}_k = \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}} (\mathbf{g}_k^i - \mathbf{h}_k^i) + \mathbf{h}_k^i \right)$.

With no memory ($\mathbf{h}_k^i = 0$ for any k in \mathbb{N}^*):

$$\mathbb{E} \|\tilde{\mathbf{g}}_k\|^2 \leq \frac{A}{N^2} \sum_{i=0}^N \mathbb{E} \|\mathbf{g}_k^i\|^2 + \frac{B}{N^2} \sum_{i=0}^N \mathbb{E} \|\mathbf{g}_k^i - \nabla F_i(w_*)\|^2 + L \langle \nabla F(w_k), w_k - w_* \rangle.$$

With memory:

$$\begin{aligned} \mathbb{E} \|\tilde{\mathbf{g}}_k\|^2 &\leq \frac{A}{N^2} \sum_{i=1}^N \mathbb{E} \|\mathbf{g}_k^i - \mathbf{g}_{k,*}^i\|^2 + \frac{B}{N^2} \sum_{i=1}^N \mathbb{E} \|\mathbf{h}_k^i - \nabla F_i(w_*)\|^2 \\ &\quad + L \langle \nabla F(w_k), w_k - w_* \rangle + \frac{C\sigma_*}{Nb}. \end{aligned}$$

We note $\tilde{g}_k = \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}} (g_k^i - h_k^i) + h_k^i \right)$.

With no memory ($h_k^i = 0$ for any k in \mathbb{N}^*):

$$\mathbb{E} \|\tilde{g}_k\|^2 \leq \frac{A}{N^2} \sum_{i=0}^N \mathbb{E} \|g_k^i\|^2 + \frac{B}{N^2} \sum_{i=0}^N \mathbb{E} \|g_k^i - \nabla F_i(w_*)\|^2 + L \langle \nabla F(w_k), w_k - w_* \rangle.$$

With memory:

$$\begin{aligned} \mathbb{E} \|\tilde{g}_k\|^2 &\leq \frac{A}{N^2} \sum_{i=1}^N \mathbb{E} \|g_k^i - g_{k,*}^i\|^2 + \frac{B}{N^2} \sum_{i=1}^N \mathbb{E} \|h_k^i - \nabla F_i(w_*)\|^2 \\ &\quad + L \langle \nabla F(w_k), w_k - w_* \rangle + \frac{C\sigma_*}{Nb}. \end{aligned}$$

- $\langle \nabla F(w_k), w_k - w_* \rangle$ allows to use strong-convexity,
- $\|g_k^i\|^2$ makes appears the constant of heterogeneity B^2 !

Backup on MCM

A practical algorithm?



Ghost cannot be implemented in practice!

⇒ Which choice do we have?

A practical algorithm?



Ghost cannot be implemented in practice!

⇒ Which choice do we have?

Ghost

$$w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$
$$\hat{w}_k = w_{k-1} - \gamma C_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$

A practical algorithm?



Ghost cannot be implemented in practice!

⇒ Which choice do we have?

Ghost

$$w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$

$$\hat{w}_k = w_{k-1} - \gamma C_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$

Update compression

$$w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$

$$\hat{w}_k = \hat{w}_{k-1} - \gamma C_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N C_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$

A practical algorithm?



Ghost cannot be implemented in practice!

⇒ Which choice do we have?

Ghost

$$w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$
$$\hat{w}_k = w_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$

Model compression

$$w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$
$$\hat{w}_k = \mathcal{C}_{\text{down}}(w_k)$$

Update compression

$$w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$
$$\hat{w}_k = \hat{w}_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$

A practical algorithm?



Ghost cannot be implemented in practice!

⇒ Which choice do we have?

Ghost

$$w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$
$$\hat{w}_k = w_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$

Update compression

$$w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$
$$\hat{w}_k = \hat{w}_{k-1} - \gamma \mathcal{C}_{\text{down}} \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$

Model compression

$$w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$
$$\hat{w}_k = \mathcal{C}_{\text{down}}(w_k)$$

Model difference compression

$$w_k = w_{k-1} - \gamma \left(\frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i(\hat{w}_{k-1})) \right)$$
$$\hat{w}_k = \hat{w}_{k-1} - \mathcal{C}_{\text{down}}(w_k - \hat{w}_{k-1})$$

First attempts - Variance of the local iterate is too high.



- Update compression
- Model difference compression
- Model compression
- MCM

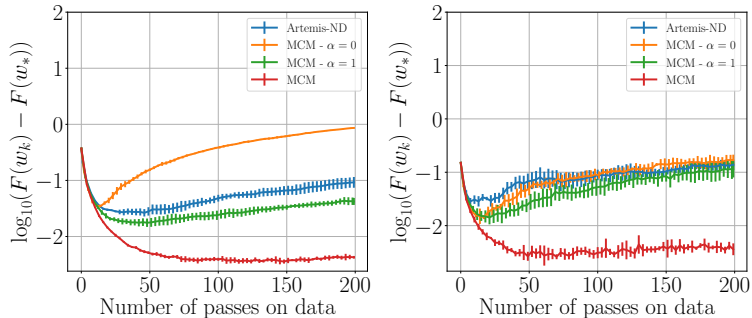


Figure 11: Comparing MCM on two datasets with three other algorithms using a non-degraded update, $\gamma = 1/L$.

Smoothed version of F :

$$F_\rho(w) := \mathbb{E}[F(w + \rho X)], \text{ with } X \sim \mathcal{N}(0, I).$$

Smoothed version of F :

$$F_\rho(w) := \mathbb{E}[F(w + \rho X)], \text{ with } X \sim \mathcal{N}(0, I).$$

$\nabla F(\hat{w}_{k-1})$ can be considered as an unbiased gradient of the smoothed function F_ρ at point w_{k-1} , with : $F_\rho : w \mapsto \mathbb{E}[F(w - w_{k-1} + \hat{w}_{k-1})]$

Smoothed version of F :

$$F_\rho(w) := \mathbb{E}[F(w + \rho X)], \text{ with } X \sim \mathcal{N}(0, I).$$

$\nabla F(\hat{w}_{k-1})$ can be considered as an unbiased gradient of the smoothed function F_ρ at point w_{k-1} , with : $F_\rho : w \mapsto \mathbb{E}[F(w - w_{k-1} + \hat{w}_{k-1})]$ i.e.:

$$\nabla F(\hat{w}_{k-1}) = \nabla F_\rho(w_{k-1})$$

Smoothed version of F :

$$F_\rho(w) := \mathbb{E}[F(w + \rho X)], \text{ with } X \sim \mathcal{N}(0, I).$$

$\nabla F(\hat{w}_{k-1})$ can be considered as an unbiased gradient of the smoothed function F_ρ at point w_{k-1} , with : $F_\rho : w \mapsto \mathbb{E}[F(w - w_{k-1} + \hat{w}_{k-1})]$ i.e.:

$$\nabla F(\hat{w}_{k-1}) = \nabla F_\rho(w_{k-1})$$

Then $\mathbb{E} \langle \nabla F(\hat{w}_{k-1}), w_{k-1} - w_* \rangle = \mathbb{E} \langle \nabla F_\rho(w_{k-1}), w_{k-1} - w_* \rangle$ which is the quantity that appears when developing the squared-norm of the update equation in the proof:

$$\mathbb{E} \|w_k - w_*\|^2 \leq \mathbb{E} \|w_{k-1} - w_*\|^2 - 2\gamma \langle \nabla F(\hat{w}_{k-1}), w_{k-1} - w_* \rangle + \gamma^2 \mathbb{E} \|\tilde{g}_k\|^2.$$

Smoothed version of F :

$$F_\rho(w) := \mathbb{E}[F(w + \rho X)], \text{ with } X \sim \mathcal{N}(0, I).$$

$\nabla F(\hat{w}_{k-1})$ can be considered as an unbiased gradient of the smoothed function F_ρ at point w_{k-1} , with $F_\rho : w \mapsto \mathbb{E}[F(w - w_{k-1} + \hat{w}_{k-1})]$ i.e.:

$$\nabla F(\hat{w}_{k-1}) = \nabla F_\rho(w_{k-1})$$

Then $\mathbb{E}\langle \nabla F(\hat{w}_{k-1}), w_{k-1} - w_* \rangle = \mathbb{E}\langle \nabla F_\rho(w_{k-1}), w_{k-1} - w_* \rangle$ which is the quantity that appears when developing the squared-norm of the update equation in the proof:

$$\mathbb{E} \|w_k - w_*\|^2 \leq \mathbb{E} \|w_{k-1} - w_*\|^2 - 2\gamma \langle \nabla F(\hat{w}_{k-1}), w_{k-1} - w_* \rangle + \gamma^2 \mathbb{E} \|\tilde{g}_k\|^2.$$

But two main differences:

- Objective function already smooth,
- Noise not Gaussian: we suffer from the noise because of compression and can not control it.

Let:

- $V_k = \mathbb{E}[\|w_k - w_*\|^2] + 32\gamma L \omega_{\text{down}}^2 \|w_k - H_{k-1}\|^2$
- $\Phi(\gamma) := (\omega_{\text{up}} + 1)(1 + 64\gamma L \omega_{\text{down}}^2)$

Theorem 7 (Convergence of MCM, convex case for any step-size γ)

Under all previous assumptions, for k in \mathbb{N}^* , for any $\gamma \leq \gamma_{\max}$, we have, for $\bar{w}_k = \frac{1}{k} \sum_{i=0}^{k-1} w_i$,

$$\begin{aligned} \gamma \mathbb{E}[F(w_{k-1}) - F(w_*)] &\leq V_{k-1} - V_k + \frac{\gamma^2 \sigma^2 \Phi(\gamma)}{Nb} \\ \implies \mathbb{E}[F(\bar{w}_k) - F_*] &\leq \frac{V_0}{\gamma k} + \frac{\gamma \sigma^2 \Phi(\gamma)}{Nb}. \end{aligned}$$

For a constant γ ,

- the variance term is upper bounded by

$$\frac{\gamma^2 \sigma^2}{Nb} (\omega_{\text{up}} + 1) (1 + 64\gamma L \omega_{\text{down}}^2).$$

- impact of the downlink compression is attenuated by a factor γ . As γ decreases, this makes the limit variance similar to the one of Diana [MGTR19], i.e. without downlink compression:

$$\frac{\gamma^2 \sigma^2}{Nb} (\omega_{\text{up}} + 1).$$

- This is much lower than the variance for previous algorithms using double compression:

$$\frac{\gamma^2 \sigma^2}{Nb} (\omega_{\text{up}} + 1) (\omega_{\text{down}} + 1).$$

Maximal learning rate to ensure convergence:

$$\gamma_{\max} := \min(\gamma_{\max}^{\text{up}}, \gamma_{\max}^{\text{down}}, \gamma_{\max}^{\text{Y}})$$

where:

1. $\gamma_{\max}^{\text{up}} := (2L(1 + \omega_{\text{up}}/N))^{-1}$ corresponds to the classical constraint on the learning rate in the unidirectional regime,
2. $\gamma_{\max}^{\text{down}} := (8L\omega_{\text{down}})^{-1}$ comes from the downlink compression,
3. $\gamma_{\max}^{\text{Y}} := (8\sqrt{2}L\omega_{\text{down}}\sqrt{8\omega_{\text{down}} + \omega_{\text{up}}/N})^{-1}$ is a combined constraint that arises when controlling the variance term $\|w_k - H_k\|^2$.

Maximal learning rate to ensure convergence:

$$\gamma_{\max} := \min(\gamma_{\max}^{\text{up}}, \gamma_{\max}^{\text{down}}, \gamma_{\max}^{\text{Y}})$$

where:

1. $\gamma_{\max}^{\text{up}} := (2L(1 + \omega_{\text{up}}/N))^{-1}$ corresponds to the classical constraint on the learning rate in the unidirectional regime,
2. $\gamma_{\max}^{\text{down}} := (8L\omega_{\text{down}})^{-1}$ comes from the downlink compression,
3. $\gamma_{\max}^{\text{Y}} := (8\sqrt{2}L\omega_{\text{down}}\sqrt{8\omega_{\text{down}} + \omega_{\text{up}}/N})^{-1}$ is a combined constraint that arises when controlling the variance term $\|w_k - H_k\|^2$.

Remarks:

- constraints are weaker than in the “degraded” framework

$$\gamma_{\max}^{\text{Dore}} \leq (8L(1 + \omega_{\text{down}})(1 + \omega_{\text{up}}/N))^{-1},$$

- if $\omega_{\text{up}, \text{down}} \rightarrow \infty$ and $\omega_{\text{down}} \simeq \omega_{\text{up}} \simeq \omega$, the maximal learning rate for MCM is $(L\omega^{3/2})^{-1}$, while it is $(L\omega^2)^{-1}$ for Dore/Artemis.

Our γ_{\max} is thus larger by a factor $\sqrt{\omega}$

Rates, complexities, and maximal step size for Diana, Artemis, Dore and MCM.

Table 3: Summary of rates on the initial condition, limit variance, asympt. complexities and γ_{\max} .

Problem		Diana	Artemis, Dore	MCM
	$L\gamma_{\max} \propto$	$1/(\omega_{\text{up}} + 1)$	$1/(\omega_{\text{up}} + 1)(\omega_{\text{down}} + 1)$	$1/(\omega_{\text{down}} + 1)\sqrt{\omega_{\text{up}} + 1} \wedge 1/(\omega_{\text{up}} + 1)$
	Lim. var. $\propto \gamma^2 \sigma^2 / n \times$	$(\omega_{\text{up}} + 1)$	$(\omega_{\text{up}} + 1)(\omega_{\text{down}} + 1)$	$(\omega_{\text{up}} + 1)(1 + \gamma L \omega_{\text{down}}^2)$
Str.-convex	Rate on init. cond. (SC)	$(1 - \gamma\mu)^k$	$(1 - \gamma\mu)^k$	$(1 - \gamma\mu)^k$
	Complexity	$(\omega_{\text{up}} + 1)/\mu\epsilon N$	$(\omega_{\text{up}} + 1)(\omega_{\text{down}} + 1)/\mu\epsilon N$	$(\omega_{\text{up}} + 1)/\mu\epsilon N$
Convex	Complexity	$(\omega_{\text{up}} + 1)/\epsilon^2$	$(\omega_{\text{up}} + 1)(\omega_{\text{down}} + 1)/\epsilon^2$	$(\omega_{\text{up}} + 1)/\epsilon^2$

⇒ Consists in performing independent compressions for each device.

Theorem 8

Theorem 4 is still valid for Rand-MCM

- Improvement in Rand-MCM: because we average gradients at several random points, reducing the impact of ω_{down} .
- Dominating term is independent of ω_{down} : we expect to reduce only the second-order term.

⇒ Consists in performing independent compressions for each device.

Theorem 8

Theorem 4 is still valid for Rand-MCM

- Improvement in Rand-MCM: because we average gradients at several random points, reducing the impact of ω_{down} .
- Dominating term is independent of ω_{down} : we expect to reduce only the second-order term.

Theorem 9 (Convergence in the quadratic case)

Under A1, A3, A7, with $\mu = 0$, if the function is quadratic, after running $K > 0$ iterations, for any $\gamma \leq \gamma_{\max}$, we have

$$\mathbb{E}[F(\bar{w}_K) - F_*] \leq \frac{V_0}{\gamma K} + \frac{\gamma \sigma^2 \Phi^{\text{Rd}}(\gamma)}{Nb},$$

with $\Phi^{\text{Rd}}(\gamma) = (1 + \omega_{\text{up}}) \left(1 + \frac{4\gamma^2 L^2 \omega_{\text{down}}}{K} \left(\frac{1}{\mathbf{C}} + \frac{\omega_{\text{up}}}{N} \right) \right)$ and $\mathbf{C} = N$ for Rand-MCM, $\mathbf{C} = 1$ for MCM.

⇒ Consists in performing independent compressions for each device.

Theorem 8

Theorem 4 is still valid for Rand-MCM

- Improvement in Rand-MCM: because we average gradients at several random points, reducing the impact of ω_{down} .
- Dominating term is independent of ω_{down} : we expect to reduce only the second-order term.

Theorem 9 (Convergence in the quadratic case)

Under A1, A3, A7, with $\mu = 0$, if the function is quadratic, after running $K > 0$ iterations, for any $\gamma \leq \gamma_{\text{max}}$, we have

$$\mathbb{E}[F(\bar{w}_K) - F_*] \leq \frac{V_0}{\gamma K} + \frac{\gamma \sigma^2 \Phi^{\text{Rd}}(\gamma)}{Nb},$$

with $\Phi^{\text{Rd}}(\gamma) = (1 + \omega_{\text{up}}) \left(1 + \frac{4\gamma^2 L^2 \omega_{\text{down}}}{K} \left(\frac{1}{\mathbf{C}} + \frac{\omega_{\text{up}}}{N}\right)\right)$ and $\mathbf{C} = N$ for Rand-MCM, $\mathbf{C} = 1$ for MCM.

- Quadratic functions: right hand term in Φ multiplied by an additional $\gamma \left(\frac{1}{\mathbf{C}} + \frac{\omega_{\text{up}}}{N}\right)$.
- Randomization: further reduces by a factor N this term.

Backup on the compressors' covariance

Impact of the compression on the additive noise covariance



The additive noise writes for any $k \in \{1, \dots, K\}$, as:

$$\xi_k^{\text{add}} \stackrel{\text{def.}}{=} \xi_k(0) \stackrel{\text{algo}}{=} \nabla F(w_*) - \mathcal{C}_k(g_k(w_*)) = -\mathcal{C}_k(\langle x_k, w_* \rangle - y_k)x_k = \mathcal{C}_k(\varepsilon_k x_k).$$

By definition: $\mathfrak{C}_{\text{ania}} := \mathbb{E}[(\xi_k^{\text{add}})^{\otimes 2}] = \mathbb{E}[\mathcal{C}(\varepsilon_k x_k)^{\otimes 2}]$. Note also that $\mathcal{C}(\varepsilon_k x_k) \stackrel{\text{a.s.}}{=} \varepsilon_k \mathcal{C}(x_k)$ for all operators under consideration. Consequently

$$\mathfrak{C}_{\text{ania}} = \mathbb{E}[\varepsilon_k^2 \mathcal{C}(x_k)^{\otimes 2}] = \sigma^2 \mathbb{E}[\mathcal{C}(x_k)^{\otimes 2}]. \quad (3)$$

We study the covariance of $\mathcal{C}(x_k)$, for x_k a random variable with second-moment H , more generically we study the covariance of $\mathcal{C}(E)$, for E a random vector with distribution p_M with second moment $\mathbb{E}[E^{\otimes 2}] = M$.

Definition 4 (Compressor' covariance on p_M)

We define the following operator \mathfrak{C} which returns the covariance of a random mechanism \mathcal{C} acting on a distribution $p_M \in \mathcal{P}_M$,

$$\mathfrak{C}: \begin{array}{ll} \mathbb{C} \times \mathcal{P}_M & \rightarrow \mathbb{R}^{d \times d} \\ (\mathcal{C}, p_M) & \rightarrow \mathbb{E}[\mathcal{C}(E)^{\otimes 2}], \end{array}$$

where $E \sim p_M$ and the expectation is over the joint randomness of \mathcal{C} and E , which are considered independent, that is $\mathbb{E}[\mathcal{C}(E)^{\otimes 2}] = \int_{\mathbb{R}^d} \mathbb{E}[\mathcal{C}(e)^{\otimes 2}] dp_M(e)$.

Algorithm 3 (Distributed compressed LMS)

At any step k in $\{1, \dots, K\}$, each clients i in $\{1, \dots, N\}$ observes an oracle $g_k^i(\cdot)$ of the gradient of their local objective function F_i and applies a random compression mechanism $\mathcal{C}_k^i(\cdot)$.

For any step-size $\gamma > 0$ and any $k \in \mathbb{N}^*$, the resulting sequence of iterates $(w_k)_{k \in \mathbb{N}}$ satisfies:

$$w_k = w_{k-1} - \gamma \frac{1}{N} \sum_{i=1}^N \mathcal{C}_k^i(g_k^i(w_{k-1})).$$

Equivalently, for $w \in \mathbb{R}^d$: $\xi_k(w) = \nabla F(w) - \frac{1}{N} \sum_{i=1}^N \mathcal{C}_k^i(g_k^i(w))$.

Algorithm 3 (Distributed compressed LMS)

At any step k in $\{1, \dots, K\}$, each clients i in $\{1, \dots, N\}$ observes an oracle $g_k^i(\cdot)$ of the gradient of their local objective function F_i and applies a random compression mechanism $\mathcal{C}_k^i(\cdot)$.

For any step-size $\gamma > 0$ and any $k \in \mathbb{N}^*$, the resulting sequence of iterates $(w_k)_{k \in \mathbb{N}}$ satisfies:

$$w_k = w_{k-1} - \gamma \frac{1}{N} \sum_{i=1}^N \mathcal{C}_k^i(g_k^i(w_{k-1})).$$

Equivalently, for $w \in \mathbb{R}^d$: $\xi_k(w) = \nabla F(w) - \frac{1}{N} \sum_{i=1}^N \mathcal{C}_k^i(g_k^i(w))$.

Two scenarios:

- Heterogeneous covariances: for i, j in $\{1, \dots, N\}$, possibly $H_i \neq H_j$ (covariate-shift),
- Heterogeneous optimal points: for i, j in $\{1, \dots, N\}$, possibly $w_*^i \neq w_*^j$ (optimal-point-shift).

Algorithm 3 (Distributed compressed LMS)

At any step k in $\{1, \dots, K\}$, each clients i in $\{1, \dots, N\}$ observes an oracle $g_k^i(\cdot)$ of the gradient of their local objective function F_i and applies a random compression mechanism $\mathcal{C}_k^i(\cdot)$.

For any step-size $\gamma > 0$ and any $k \in \mathbb{N}^*$, the resulting sequence of iterates $(w_k)_{k \in \mathbb{N}}$ satisfies:

$$w_k = w_{k-1} - \gamma \frac{1}{N} \sum_{i=1}^N \mathcal{C}_k^i(g_k^i(w_{k-1})).$$

Equivalently, for $w \in \mathbb{R}^d$: $\xi_k(w) = \nabla F(w) - \frac{1}{N} \sum_{i=1}^N \mathcal{C}_k^i(g_k^i(w))$.

Two scenarios:

- Heterogeneous covariances: for i, j in $\{1, \dots, N\}$, possibly $H_i \neq H_j$ (covariate-shift),
- Heterogeneous optimal points: for i, j in $\{1, \dots, N\}$, possibly $w_*^i \neq w_*^j$ (optimal-point-shift).

Corollary 1 (covariate-shift)

Theorem 6 holds.

How to compute the ania's covariance using the compressor's covariance?

We have for any clients $i, j \in \{1, \dots, N\}$, $w_*^i = w_*^j$, thus

$$\begin{aligned} \xi_k^{\text{add}} &\stackrel{\text{def. 2}}{=} \xi_k(0) \stackrel{\text{algo 3}}{=} \nabla F(w_*) - \frac{1}{N} \sum_{i=1}^N \mathbf{C}_k^i(g_k^i(w_*)) \\ &= -\frac{1}{N} \sum_{i=1}^N \mathbf{C}_k^i(\langle x_k^i, w_* \rangle - y_k^i) x_k^i \stackrel{w_*^i = w_*^j}{=} \frac{1}{N} \sum_{i=1}^N \mathbf{C}_k^i(\varepsilon_k^i x_k^i). \end{aligned}$$

Next for all operators under consideration we have $\mathbf{C}_k^i(\varepsilon_k^i x_k^i) \stackrel{\text{a.s.}}{=} \varepsilon_k^i \mathbf{C}_k^i(x_k^i)$, thus, with p_{H_i} denoting the distribution of x_k^i with covariance H_i , we have:

$$\begin{aligned} \mathfrak{C}_{\text{ania}} &= \mathbb{E}[(\xi_k^{\text{add}})^{\otimes 2}] = \mathbb{E}\left[\left(\frac{1}{N} \sum_{i=1}^N \mathbf{C}_k^i(\varepsilon_k^i x_k^i)\right)^{\otimes 2}\right] \stackrel{\text{indep. of } (\mathbf{C}_k^i)_{i=1}^d}{=} \frac{1}{N^2} \sum_{i=1}^N \mathbb{E}[\mathbf{C}_k^i(\varepsilon_k^i x_k^i)^{\otimes 2}] \\ &= \frac{\sigma^2}{N^2} \sum_{i=1}^N \mathbb{E}[\mathbf{C}_k^i(x_k^i)^{\otimes 2}] \stackrel{\text{Def. 4}}{=} \frac{\sigma^2}{N^2} \sum_{i=1}^N \mathfrak{C}(\mathbf{C}_k^i, p_{H_i}) \stackrel{\text{notation}}{=} \frac{\sigma^2}{N} \overline{\mathfrak{C}((\mathbf{C}_k^i, p_{H_i})_{i=1}^N)}. \end{aligned} \quad (4)$$

The operator $\overline{\mathfrak{C}((\mathbf{C}_k^i, p_{H_i})_{i=1}^N)}$ generalizes the notion of *compressor's covariance* (Definition 4).

By definition, we have:

$$\xi_k(w - w_*) \stackrel{\text{Def. 1\&Alg.3}}{=} H_F(w - w_*) - \frac{1}{N} \sum_{i=1}^N \mathbf{c}^i(g_k^i(w)), \text{ thus } \xi_k^{\text{add}} \stackrel{\text{Def. 2}}{=} \frac{1}{N} \sum_{i=1}^N \mathbf{c}^i(g_{k,*}^i),$$

with $g_{k,*}^i = (x_k^i \otimes x_k^i)(w_* - w_*^i) + x_k^i \varepsilon_k^i$. We thus have, for any $k \in \mathbb{N}$:

$$\begin{aligned} \mathfrak{C}_{\text{ania}} &= \mathbb{E}[(\xi_k^{\text{add}})^{\otimes 2}] \stackrel{\nabla F(w_*)=0}{=} \mathbb{E}\left[\left(\frac{1}{N} \sum_{i=1}^N \mathbf{c}^i(g_{k,*}^i) - \nabla F_i(w_*)\right)^{\otimes 2}\right] \\ &\stackrel{\substack{\forall i \neq j, \mathbf{c}_k^i \perp \mathbf{c}_k^j \\ \mathbb{E} \mathbf{c}_k^i(g_{k,*}^i) = \nabla F_i(w_*)}}{=} \frac{1}{N^2} \sum_{i=1}^N \mathbb{E}\left[(\mathbf{c}_k^i(g_{k,*}^i) - \nabla F_i(w_*))^{\otimes 2}\right] \\ &= \frac{1}{N^2} \sum_{i=1}^N (\mathbb{E}[\mathbf{c}_k^i(g_{k,*}^i)^{\otimes 2}] - \nabla F_i(w_*)^{\otimes 2}) \\ &= \frac{\sigma^2}{N^2} \sum_{i=1}^N \mathfrak{C}(\mathbf{c}^i, p_{\Theta_i}) - \frac{1}{N^2} H \sum_{i=1}^N (w_* - w_*^i)^{\otimes 2} H \leq \frac{\sigma^2}{N} \mathfrak{C}((\mathbf{c}^i, p_{\Theta_i})_{i=1}^N), \end{aligned}$$

where p_{Θ_i} is the distribution of $g_{k,*}^i$ (for any k).

In order to bound this quantity, following [DFB17], we make the following assumption.

Assumption 8

The kurtosis for the projection of the covariates x_1^i (or equivalently x_k^i for any k) is bounded on any direction $z \in \mathbb{R}^d$, i.e., there exists $\kappa > 0$, such that:

$$\forall i \in \{1, \dots, N\}, \forall z \in \mathbb{R}^d, \quad \mathbb{E} \left[\langle z, x_1^i \rangle^4 \right] \leq \kappa \langle z, Hz \rangle^2$$

Proposition 1 (Impact of client-heterogeneity.)

Let W_* be a random variable uniformly distributed over $\{w_*^i, i \in \{1, \dots, N\}\}$, thus such that, $\text{Cov}[W_*] = \frac{1}{N} \sum_{i=1}^N (w_* - w_*^i)^{\otimes 2}$, then:

$$\frac{1}{N} \sum_{i=1}^N \Theta_i \leq (\kappa \text{Tr}(H \text{Cov}[W_*]) + \sigma^2) H.$$

1) Before compression is possibly applied, the noise remains structured, i.e., with covariance proportional to H , in the case of concept-shift

2) Compared to the homogeneous case, the averaged second-order moment increases from $\sigma^2 H$ to $(\kappa \text{Tr}(H \text{Cov}[W_*]) + \sigma^2) H$.

\implies shows impact of the dispersion of the optimal points. $(w_*^i)_{i=1}^N$.

Heterogeneous optimal points w_*^i with memory



Artemis with only uplink compression:

$$w_k = w_{k-1} - \gamma \frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i - h_k^i) + h_k^i$$

$$h_{k+1}^i = h_k^i + \alpha \mathcal{C}_{\text{up}}(g_k^i - h_k^i),$$

Artemis with only uplink compression:

$$w_k = w_{k-1} - \gamma \frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i - h_k^i) + h_k^i$$

$$h_{k+1}^i = h_k^i + \alpha \mathcal{C}_{\text{up}}(g_k^i - h_k^i),$$

⚠ Random fields are no more i.i.d. \implies Definition 1 is no more fulfilled, invalidating Theorem 6. ⚠

Artemis with only uplink compression:

$$w_k = w_{k-1} - \gamma \frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i - h_k^i) + h_k^i$$
$$h_{k+1}^i = h_k^i + \alpha \mathcal{C}_{\text{up}}(g_k^i - h_k^i),$$

⚠ Random fields are no more i.i.d. \implies Definition 1 is no more fulfilled, invalidating Theorem 6. ⚠

Theorem 10 (CLT for concept-shift heterogeneity)

Under some assumption, with $\mu > 0$, for any step-size $(\gamma_k)_{k \in \mathbb{N}^*}$ s.t. $\gamma_k = 1/\sqrt{k}$. Then

1. $(\sqrt{K} \bar{\eta}_{K-1})_{K>0} \xrightarrow{K \rightarrow +\infty} \mathcal{L} \mathcal{N}(0, H_F^{-1} \mathfrak{C}_{\text{ania}}^\infty H_F^{-1})$,
2. $\mathfrak{C}_{\text{ania}}^\infty = \overline{\mathfrak{C}((C^i, p_{\Theta_i'})_{i=1}^N)}$, where $p_{\Theta_i'}$ is the distribution of $g_{k,*}^i - \nabla F_i(w_*)$.

Artemis with only uplink compression:

$$w_k = w_{k-1} - \gamma \frac{1}{N} \sum_{i=1}^N \mathcal{C}_{\text{up}}(g_k^i - h_k^i) + h_k^i$$

$$h_{k+1}^i = h_k^i + \alpha \mathcal{C}_{\text{up}}(g_k^i - h_k^i),$$

⚠ Random fields are no more i.i.d. \implies Definition 1 is no more fulfilled, invalidating Theorem 6. ⚠

Theorem 10 (CLT for concept-shift heterogeneity)

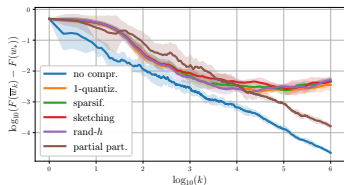
Under some assumption, with $\mu > 0$, for any step-size $(\gamma_k)_{k \in \mathbb{N}^*}$ s.t. $\gamma_k = 1/\sqrt{k}$. Then

1. $(\sqrt{K} \bar{\eta}_{K-1})_{K > 0} \xrightarrow{K \rightarrow +\infty} \mathcal{L} \rightarrow \mathcal{N}(0, H_F^{-1} \mathfrak{C}_{\text{ania}}^\infty H_F^{-1}),$
2. $\mathfrak{C}_{\text{ania}}^\infty = \overline{\mathfrak{C}((C^i, p_{\Theta_i'}^N)_{i=1}^N)},$ where $p_{\Theta_i'}$ is the distribution of $g_{k,*}^i - \nabla F_i(w_*)$.

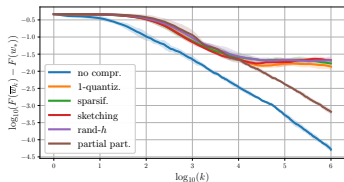
1. Settings of heterogeneous optimal points $(w_*^i)_{i=1}^N$: convergence still impacted by heterogeneity but with smaller additive noise's covariance as $\Theta_i' < \Theta_i$.
2. Deterministic gradients (batch case), we case $\Theta_i' \equiv 0$.
3. Recover asymptotically the results stated by Theorem 6 in the general setting of i.i.d. random fields $(\xi_k(\eta_{k-1}))_{k \in \mathbb{N}^*}$.

Covariate-shift

Synthetic dataset

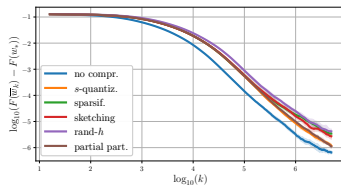


(a) No shift: $\forall i \in \{1, \dots, N\}, H_i = H$

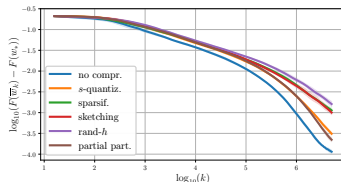


(d) $\forall i, j \in \{1, \dots, N\}, H_i \neq H_j$

Real datasets



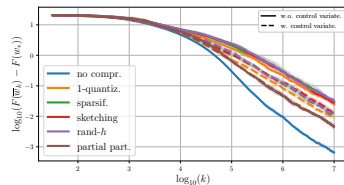
(b) quantum



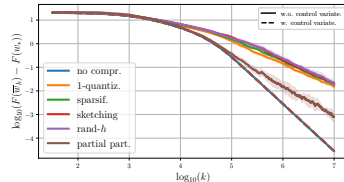
(e) cifar-10

Concept-shift

Synthetic dataset



(c) Batch stochastic gradient



(f) True gradient $g_k^i = \nabla F_i$

Figure 12: Logarithm excess loss of the Polyak-Ruppert iterate iterations for $N = 10$ clients.