
Generalization under Byzantine & Poisoning Attacks: Tight Stability Bounds in Robust Distributed Learning

Thomas Boudou

Inria, Université de Montpellier, INSERM
Montpellier, France
thomas.boudou@inria.fr

Batiste Le Bars

Inria, Université de Lille
Lille, France
batiste.le-bars@inria.fr

Nirupam Gupta

University of Copenhagen
Copenhagen, Denmark
nigu@di.ku.dk

Aurélien Bellet

Inria, Université de Montpellier, INSERM
Montpellier, France
aurelien.bellet@inria.fr

Abstract

Robust distributed learning algorithms aim to maintain good performance in distributed and federated settings, even in the presence of misbehaving workers. Two primary threat models have been studied: *Byzantine attacks*, where misbehaving workers can send arbitrarily corrupted updates, and *data poisoning attacks*, where misbehavior is limited to manipulation of local training data. While prior work has shown comparable optimization error under both threat models, a fundamental question remains open: *How do these threat models impact generalization?* Empirical evidence suggests a gap between the two threat models, yet it remains unclear whether it is fundamental or merely an artifact of suboptimal attacks. In this work, we present the first theoretical investigation into this problem, formally showing that Byzantine attacks are intrinsically more harmful to generalization than data poisoning. Specifically, we prove that: (i) under data poisoning, the uniform algorithmic stability of a robust distributed learning algorithm, with optimal optimization error, degrades by an additive factor of $\Theta(\frac{f}{n-f})$, with f the number of misbehaving workers out of n ; and (ii) In contrast, under Byzantine attacks, the degradation is in $\mathcal{O}(\sqrt{\frac{f}{n-2f}})$. This difference in stability leads to a generalization error gap that is especially significant as f approaches its maximum value $\frac{n}{2}$.

1 Introduction

With the proliferation of large-scale data and pervasive connectivity, distributed learning has evolved from a theoretical concept into a fundamental paradigm in modern AI [24]. From federated smartphones to geo-distributed datacenters, algorithms now orchestrate millions of workers toward shared objectives. Yet this democratization of learning comes at a cost: unlike the idealized setting of perfectly cooperative workers, real-world systems face a range of adversarial behaviors, broadly categorized under the umbrella of Byzantine failures [20]. These failures, encompassing everything from intermittent data corruption to sophisticated, coordinated attacks, pose a significant threat to the convergence, accuracy, and overall stability of distributed learning algorithms. Robust distributed learning [21] aims to maintain strong learning guarantees in the presence of such failures.

In the Byzantine failure model, a misbehaving worker can send arbitrarily corrupted model updates to the server, while in data poisoning attacks, its influence is limited to manipulating local training data. Although prior work formally shows that both threat models cause similar degradation in

empirical risk [15], experiments consistently find that the generalization error—*i.e.*, expected risk on unseen data—is significantly worse under Byzantine failures [4, 26]. It remains unclear from existing literature whether this gap reflects a fundamental difference between the two threat models or the suboptimality of practical attacks. We address this open question with the first theoretical analysis of the generalization gap between Byzantine and data poisoning attacks.

To this end, we use algorithmic stability [10, 13] as our analytical tool for quantifying resilience to both forms of attacks. Algorithmic stability, which measures an algorithm’s sensitivity to changes in a single training example, offers key advantages over alternatives techniques for characterizing generalization. Unlike uniform convergence bounds [42, 7], which typically rely on i.i.d. assumptions and restricted hypothesis spaces, or PAC-Bayes [32, 6] and mutual information bounds [34], which provide less direct insight into algorithm-specific behavior, algorithmic stability takes an algorithm-centric view of generalization. This makes it particularly well-suited for robust distributed learning, where it enables a focused analysis of how specific attacks and robust aggregation rules influence the generalization performance of the overall learning procedure.

Our contributions. This paper provides, to our knowledge, the first theoretically grounded explanation for the gap in generalization performance between two threat models in robust distributed learning: Byzantine and data poisoning attacks. Comparisons with related work are deferred to Section 5. Our main contributions—beyond formalizing a unified framework (Section 2)—are:

- (i) We derive general upper bounds on uniform stability for Byzantine attacks under convex and strongly convex objectives (Section 3.1). Specifically, for a convex loss function, our analysis shows that the uniform algorithmic stability of robust distributed GD and SGD incorporating SMEA [5]—an optimal high-dimensional aggregation rule— under Byzantine attacks is in $\mathcal{O}(\sqrt{\frac{f}{n-2f}})$, f being the number of misbehaving workers out of n .
- (ii) We establish matching upper and lower bounds on the uniform algorithmic stability of this algorithm under data poisoning (Section 3.2). Specifically, we show that the bound improves to $\Theta(\frac{f}{n-f})$. Since generalization error is bounded by stability [10], this suggests that Byzantine attacks can be substantially more harmful to generalization than data poisoning, particularly when f is close to $n/2$, the theoretical limit on the number of misbehaving workers.
- (iii) For smooth nonconvex loss functions, this discrepancy persists in the upper bounds. Under poisoning attacks, we derive sharper stability bounds for robust distributed SGD. This improvement—particularly in its dependence on the number of local samples—is enabled by aggregation rules that preserve the smoothness of the loss function (Section 3.3).
- (iv) We numerically validate (Section 4) the tightness of our upper bound under Byzantine attacks in the convex case, by designing a tailored Byzantine attack that yields higher empirical instability than an optimal poisoning attack. Remarkably, this translates directly into a generalization gap between the two attacks. This reveals a fundamental difference: Byzantine failures can cause strictly worse generalization than data poisoning—a gap that persists even when both attacks are optimally crafted.

Practical implications of our results. Our theoretical results highlight the practical relevance of cryptographic tools—particularly *zero-knowledge proofs* (ZKPs) [18, 40]—in the context of distributed learning. ZKPs can prevent Byzantine failures by verifying that each worker’s input data lies within an appropriate range and that their model updates remain consistent with their initially committed dataset, all without revealing the data itself [35, 1, 37]. However, they do not guard against data poisoning, which stems from malicious or corrupted local data. Therefore, cryptographic safeguards must be paired with robust learning algorithms that can tolerate compromised training data to ensure end-to-end resilience.

2 Preliminaries

2.1 Problem setting

We consider a distributed setup with a central server and n workers, where up to f (with $f < \frac{n}{2}$) may be subject to either *Byzantine failure* or *data poisoning* attacks, as defined below. We refer to such workers as *misbehaving*, unless we need to distinguish between the two attack types. Their

identities are unknown to the server. While the actual number of misbehaving workers may be less than f , we assume the worst-case scenario where exactly f workers are misbehaving. The remaining *honest* workers form the set $\mathcal{H} \subseteq [n] := \{1, \dots, n\}$, with $|\mathcal{H}| = n - f$. Each honest worker $i \in \mathcal{H}$ holds a local dataset $\mathcal{D}_i = \{z^{(i,1)}, \dots, z^{(i,m)}\}$ composed of m i.i.d. data points (or samples) from an input space \mathcal{Z} drawn from a distribution p_i . We denote $\mathcal{S} = \cup_{i \in \mathcal{H}} \mathcal{D}_i$. Given a parameter vector $\theta \in \Theta \subset \mathbb{R}^d$ representing the model, a data point $z \in \mathcal{Z}$ incurs a differentiable loss defined by a real-valued function $\ell(\theta; z)$. The goal is to minimize the *population* risk over the honest workers:

$$R_{\mathcal{H}}(\theta) = \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \mathbb{E}_{z \sim p_i} [\ell(\theta, z)].$$

As we only have access to a finite number of samples from each distribution, the above objective can only be solved approximately by minimizing the *empirical* risk over the honest workers:

$$\hat{R}_{\mathcal{H}}(\theta) = \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \hat{R}_i(\theta), \quad \text{with} \quad \hat{R}_i(\theta) = \frac{1}{m} \sum_{z \in \mathcal{D}_i} \ell(\theta, z).$$

The population risk $R_{\mathcal{H}}(\theta)$ of a distributed learning algorithm's output $\mathcal{A}(\mathcal{S})$ can be decomposed into the generalization and optimization errors (see, e.g. [9, 23]), i.e.,

$$\mathbb{E}[R_{\mathcal{H}}(\mathcal{A}(\mathcal{S})) - \inf_{\theta \in \Theta} R_{\mathcal{H}}(\theta)] \leq \mathbb{E}[R_{\mathcal{H}}(\mathcal{A}(\mathcal{S})) - \hat{R}_{\mathcal{H}}(\mathcal{A}(\mathcal{S}))] + \mathbb{E}[\hat{R}_{\mathcal{H}}(\mathcal{A}(\mathcal{S})) - \inf_{\theta \in \Theta} \hat{R}_{\mathcal{H}}(\theta)] \quad (1)$$

Prior work primarily only considers the optimization error, i.e., the second term on the right-hand side. In contrast, we are interested in minimizing the honest population risk under Byzantine attacks. Formally, we define the notion of *robustness* as follows. We denote the Euclidean norm by $\|\cdot\|_2$.

Definition 2.1 ($(f, \rho, \{\text{empirical, statistical}\})$ -resilience). *A distributed algorithm \mathcal{A} is said to be $(f, \rho, \text{statistical})$ -resilient if it enables the server to output a parameter vector $\hat{\theta}$ such that $\mathbb{E}_{\mathcal{A}}[R_{\mathcal{H}}(\hat{\theta}) - \inf_{\theta \in \Theta} R_{\mathcal{H}}(\theta)] \leq \rho$. In the case of nonconvex population risk, we aim for $\mathbb{E}_{\mathcal{A}}\|\nabla R_{\mathcal{H}}(\hat{\theta})\|_2^2 \leq \rho$. We say $(f, \rho, \text{empirical})$ -resilient when we instead consider empirical risk.*

Note that the (standard) notion of $(f, \rho, \text{empirical})$ -resilience—see, e.g., [30]—is impossible in general (for any ρ) when $f \geq n/2$, justifying our assumption $f < n/2$.

Threat models. We aim to analyze this resilience property under two standard threat models widely studied in the literature [14, 38].

Byzantine failure attacks. In this threat model, misbehaving workers can act arbitrarily: they may deviate from the prescribed algorithm and send adversarial updates to the server. These Byzantine workers are omniscient, with access to all information exchanged between honest workers and the server. This model follows the classical notion of Byzantine failures in distributed systems introduced in the seminal work of [27].

Data poisoning attacks. In this restricted threat model, misbehaving workers follow the prescribed algorithm but may arbitrarily corrupt their local training data *before* execution. Specifically, for each misbehaving worker $i \notin \mathcal{H}$, a local dataset \mathcal{D}_i is sampled from a distribution q_i over \mathcal{Z} , where q_i is adversarially chosen with full knowledge of the honest workers' data distributions $\{p_i, i \in \mathcal{H}\}$.

2.2 Background on robust distributed optimization

Robust distributed optimization algorithms aim to achieve $(f, \rho, \text{empirical})$ -resilience. They are typically adaptations of standard first-order iterative optimization algorithms like gradient descent (GD) or stochastic gradient descent (SGD). These algorithms are made robust by replacing the server-side averaging operator by a robust aggregation rule $F : (\mathbb{R}^d)^n \rightarrow \mathbb{R}^d$ [20]. Formally, given a learning rate γ , the parameter θ_t at iteration $t \in \{0, \dots, T-1\}$ is updated as follows:

$$\theta_{t+1} = G_{\gamma}^F(\theta_t) := \theta_t - \gamma F(g_t^{(1)}, \dots, g_t^{(n)}),$$

where $g_t^{(i)}$ denotes the update sent by worker i at iteration t . For an honest worker $i \in \mathcal{H}$, we have $g_t^{(i)} = \nabla \hat{R}_i(\theta_t)$ under GD, or $g_t^{(i)} = \nabla \ell(\theta_t; z_t^{(i)})$ under SGD, with $z_t^{(i)}$ uniformly sampled from the local dataset \mathcal{D}_i . For a misbehaving worker $i \notin \mathcal{H}$, $g_t^{(i)}$ is either an arbitrary vector in \mathbb{R}^d (Byzantine failure) or a gradient computed on corrupted data (poisoning attack).

Recent advances in robust distributed optimization have identified key properties that a robust aggregation rule must satisfy to ensure $(f, \rho, \text{empirical})$ -resilience [4, 5]. The following definition encompasses a wide range of aggregation schemes and has been shown to yield tight resilience guarantees across a variety of practical distributed learning settings.

Definition 2.2 ($(f, \kappa, \|\cdot\|_{op})$ -robustness). Let $\|\cdot\|_{op}$ denote an operator norm for p.s.d. matrices. Let $n \geq 1$, $0 \leq f < n/2$ and $\kappa \geq 0$. An aggregation rule F is $(f, \kappa, \|\cdot\|_{op})$ -robust if for any $g_1, \dots, g_n \in \mathbb{R}^d$, any set $S \subset [n]$ of size $n - f$, with $\bar{g}_S = \frac{1}{|S|} \sum_{i \in S} g_i$ and $\Sigma_S = \frac{1}{|S|} \sum_{i \in S} (g_i - \bar{g}_S)(g_i - \bar{g}_S)^\top$,

$$\|F(g_1, \dots, g_n) - \bar{g}_S\|_2^2 \leq \kappa \|\Sigma_S\|_{op}.$$

Depending on the context, we either use the nuclear norm (i.e., the trace Tr) or the spectral norm $\|\cdot\|_{sp}$ (i.e., the largest eigenvalue) for $\|\cdot\|_{op}$. f and κ are referred to as the robustness parameter and robustness coefficient of F , respectively.

An aggregation rule F that is $(f, \mathcal{O}(f/n), \|\cdot\|_{sp})$ -robust—such as *smallest maximum eigenvalue averaging* (SMEA, which averages the $n - f$ gradients that are most directionally consistent, discarding outliers based on covariance; see Definition C.1)—achieves optimal $(f, \rho, \text{empirical})$ -resilience. Specifically, iterative methods using such F attain an optimization error that matches the information-theoretic lower bound for first-order optimization under f Byzantine workers [5]. When the empirical covariance matrix of honest workers’ gradients has uniformly bounded trace over Θ , tight optimization error guarantees can also be obtained under the weaker $(f, \mathcal{O}(f/n), \text{Tr})$ -robust property. This weaker condition can be satisfied by computationally cheaper rules, such as *coordinate-wise trimmed mean* (CWTM, which removes the largest and smallest f values in each coordinate and averages the remaining entries; see Definition C.2) [4]. As noted in our comparison to related work (Section 5), prior analyses focus on optimization error, failing to account for discrepancies in generalization performance under different threat models. We address this limitation by studying how robust aggregation impacts the stability of Byzantine-resilient distributed learning algorithms.

2.3 Stability under misbehaving workers

The main objective of this work is to analyze the generalization error (eq. (1)) of robust distributed optimization algorithms. We do so through the lens of uniform algorithmic stability, which measures how sensitive a learning algorithm is to changes in its training data. This approach originates in early works [41, 12, 33] and was later popularized in [10, 36, 23]. We revisit this framework to study generalization in robust distributed learning, where the change is only in the honest workers’ data.

Definition 2.3 (Stability with misbehaving workers). Consider n workers, among which f are misbehaving, and $n - f$ are honest with local datasets of size m . A distributed randomized algorithm \mathcal{A} is said to be ε -uniformly stable if for all pairs of honest datasets $\mathcal{S}, \mathcal{S}' \in \mathcal{Z}^{(n-f)m}$ that differ in at most one sample (referred to as neighboring datasets), we have $\sup_{z \in \mathcal{Z}} \mathbb{E}[\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)] \leq \varepsilon$, where the expectation is over the randomness of \mathcal{A} .

This property leads to the following bound on the generalization error, established in Appendix A.

Proposition 2.1. If \mathcal{A} is ε -uniformly stable, then $|\mathbb{E}_{\mathcal{S}, \mathcal{A}}[R_{\mathcal{H}}(\mathcal{A}(\mathcal{S})) - \hat{R}_{\mathcal{H}}(\mathcal{A}(\mathcal{S}))]| \leq \varepsilon$.

We can now state a relationship between empirical resilience and statistical resilience: any algorithm that is $(f, \rho, \text{empirical})$ -resilient and ε -uniformly stable is $(f, \rho + \varepsilon, \text{statistical})$ -resilient. The term ε captures the impact of misbehaving workers on the generalization, and is the focus of this paper.

3 (In)stability of robust distributed learning

In this section, we first establish general upper bounds on the algorithmic stability of robust distributed GD and SGD using $(f, \kappa, \|\cdot\|_{op})$ -robust aggregation under convex and strongly convex loss functions. We then derive a sharper bound for the case of data poisoning attacks, and show that it is tight for SMEA. Finally, we provide upper bounds in the smooth, nonconvex setting, with significantly sharper results under data poisoning for CWTM—an aggregation rule that retains the loss function’s smoothness. Definitions of the relevant function regularity conditions are given in Appendix B.

3.1 Generic upper bounds under Byzantine failure attacks in convex optimization

Below, we present uniform stability upper bounds for convex optimization problems. Proofs and additional results are deferred to Appendix D.2.

Theorem 3.1 (Stability under Byzantine failure attacks in convex regimes). *Consider the setting described in Section 2 under Byzantine failure attacks. Let $\mathcal{A} \in \{\text{GD}, \text{SGD}\}$ with a $(f, \kappa, \|\cdot\|_{op})$ -robust aggregation rule F where $\|\cdot\|_{op} \in \{\text{Tr}, \|\cdot\|_{sp}\}$. Suppose $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ C -Lipschitz and L -smooth, and \mathcal{A} is run for $T \in \mathbb{N}^*$ iterations with $\gamma \leq \frac{1}{L}$.*

(i) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ convex, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq 2\gamma C^2 T \left(\frac{1}{(n-f)m} + \sqrt{\kappa} \right). \quad (2)$$

(ii) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ μ -strongly convex, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq \frac{2C^2}{\mu} \left(\frac{1}{(n-f)m} + \sqrt{\kappa} \right). \quad (3)$$

When there are no Byzantine workers ($f = 0$, hence $\kappa = 0$), our results reduce to the classical stability bounds established in [23]. However, when $f > 0$, the analysis reveals a degradation by an additive term of order $\mathcal{O}(\sqrt{\kappa})$. In the convex setting, this degradation accumulates over the T iterations, whereas in the strongly convex case, it appears as a one-time term independent of T .

Discussion on the optimality of the stability degradation. Since $\kappa \geq \frac{f}{n-2f}$ [4, Proposition 6], the additive term in Theorem 3.1 is lower bounded by $\sqrt{\frac{f}{n-2f}}$. We argue that it is not merely an artifact of our specific proof technique. In fact, for an $(f, \rho, \text{empirical})$ -resilient algorithm with a C -Lipschitz, L -smooth and μ -strongly convex loss, applying the general uniform stability result in [11, Theorem 7] yields a uniform stability upper bound of $2C\sqrt{\frac{2\rho}{\mu}} + \frac{2C^2}{\mu(n-f)m}$. Since $\rho \in \Omega(\frac{f}{n-2f}C^2)$ for this class of loss functions [4, Proposition 1], this recovers a particular case of the bound eq. (3) of Theorem 3.1 up to a constant factor. Later, in Section 4, we present numerical experiments establishing the tightness of our upper bound.

3.2 Matching lower and upper bounds under poisoning attacks in convex optimization

In this section, we derive refined stability bounds for the data poisoning threat model. In contrast to our earlier upper bounds—which accommodate the inherent unpredictability of Byzantine vectors and thus cannot leverage the loss function’s regularity—our new analysis exploits these structural properties to obtain tighter guarantees. In particular, regularity conditions such as the smoothness inequality (eq. (8) in Appendix B) for smooth nonconvex settings, and the co-coercivity inequality (eq. (9) in Appendix B) for convex and smooth problems, play a central role in stability analysis.

However, in robust distributed learning, the extent to which this regularity benefits the stability analysis depends on the specific robust aggregation rule—only the properties preserved by the rule can be effectively leveraged. Below, we present tight bounds using the SMEA aggregation rule (see Definition C.1). The focus on SMEA avoids technical complications that arise when manipulating other aggregation rules (e.g., linear loss functions cannot be used to prove a lower bound with CWTM in convex setting, see Remark D.2). However, we argue — via the *proof sketch* below — that the core reasoning behind our analysis is broadly applicable for lower bounding uniform stability under general $(f, \kappa, \|\cdot\|_{op})$ aggregation rules that fail to preserve the loss’s function co-coercivity, such as CWTM (see Remark D.3). Proofs and additional results are provided in Appendices D.4 to D.6.

We begin by stating our stability upper bounds for the data poisoning threat model, which improve upon the corresponding results for the Byzantine case presented in Theorem 3.1.

Theorem 3.2 (Stability under data poisoning attacks in convex regimes - SMEA). *Consider the setting described in Section 2 under data poisoning attacks. Let $\mathcal{A} \in \{\text{GD}, \text{SGD}\}$, with SMEA. Suppose $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ C -Lipschitz and L -smooth, and \mathcal{A} is run for $T \in \mathbb{N}^*$ iterations with $\gamma \leq \frac{1}{L}$.*

(i) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ convex, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq 2\gamma C^2 T \left(\frac{f}{n-f} + \frac{1}{(n-f)m} \right). \quad (4)$$

(ii) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ μ -strongly convex, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq \frac{2C^2}{\mu} \left(\frac{f}{n-f} + \frac{1}{(n-f)m} \right).$$

Crucially, Theorem 3.2 improves the dependence of the added instability term from the previously established $\mathcal{O}(\sqrt{\kappa})$, which evaluates at $\mathcal{O}(2\sqrt{\frac{f}{n-f}}(1 + \frac{f}{n-2f}))$ for SMEA, to a smaller $\mathcal{O}(\frac{f}{n-f})$ scaling. This refined stability bound forms the foundation of our theoretical explanation for the generalization gap between Byzantine failure and poisoning attacks, as we will see in Section 4.

We establish tightness by deriving a matching lower bound. Before stating it, we note that the lower bound is established for projected-SGD, which applies a Euclidean projection onto the positive half parameter space after each iteration. While this simplifies the lower bound proof (Appendix D.6), it does not affect the tightness of our result since the upper bound on uniform stability remains valid for projected-SGD as the projection does not increase the distance between projected points.

Theorem 3.3 (Lower bound on stability for data poisoning attacks in convex regimes). *Consider the setting described in Section 2 under data poisoning attacks. Let $\mathcal{F}_{\mathcal{Z}}$ and $\mathcal{F}_{\mathcal{Z}}^{\mu}$ denote the set of functions that are C -Lipschitz, L -smooth and convex or μ -strongly convex on \mathcal{Z} , respectively. Let $\mathcal{A} \in \{\text{GD}, \text{projected-SGD}\}$, with the SMEA. Suppose \mathcal{A} is run for $T \in \mathbb{N}^*$ iterations with $\gamma \leq \frac{1}{L}$. Then there exists a loss function $\ell \in \mathcal{F}_{\mathcal{Z}}$ and a pair of neighboring datasets such that the uniform stability of \mathcal{A} is lower bounded by*

$$\Omega\left(\gamma C^2 T \left(\frac{f}{n-f} + \frac{1}{(n-f)m}\right)\right).$$

We note that, for $\mathcal{A} = \text{projected-SGD}$, we additionally require that there exists a constant $\tau \geq c > 0$ for an arbitrary c , such that $T \geq \tau m$.

Also, for $\mathcal{A} = \text{GD}$, if $T \geq \frac{\ln(1-c)}{\ln(1-\gamma\mu)}$ for some constant $0 < c < 1$, then there exists a loss function $\ell \in \mathcal{F}_{\mathcal{Z}}^{\mu}$ and a pair of neighboring datasets such that the uniform stability of \mathcal{A} is lower bounded by

$$\Omega\left(\frac{C^2}{\mu} \left(\frac{f}{n-f} + \frac{1}{(n-f)m}\right)\right).$$

Proof sketch. We provide a general rationale for aggregation rules that fail to preserve the underlying loss's co-coercive inequality. First we provide a formal argument why an additive term may appear in any update. Let F be a robust aggregation rule and ℓ a smooth and convex loss function. Suppose there exist $\theta, \omega \in \Theta$, $\theta \neq \omega$ such that F breaks the co-coercive inequality (eq. (9)), that is:

$$\langle \theta - \omega, F(\nabla \hat{R}_1(\theta), \dots, \nabla \hat{R}_n(\theta)) - F(\nabla \hat{R}_1(\omega), \dots, \nabla \hat{R}_n(\omega)) \rangle \leq 0. \quad (5)$$

Denoting $F_{\theta} = F(\nabla \hat{R}_1(\theta), \dots, \nabla \hat{R}_n(\theta))$ and $F_{\omega} = F(\nabla \hat{R}_1(\omega), \dots, \nabla \hat{R}_n(\omega))$, eq. (5) implies:

$$\|G_{\gamma}^F(\theta) - G_{\gamma}^F(\omega)\|_2^2 = \|\theta - \omega\|_2^2 + \gamma^2 \|F_{\theta} - F_{\omega}\|_2^2 - 2\gamma \langle \theta - \omega, F_{\theta} - F_{\omega} \rangle > \|\theta - \omega\|_2^2,$$

Specifically, the quantity $\sigma_{\theta, \omega}^F = \|G_{\gamma}^F(\theta) - G_{\gamma}^F(\omega)\|_2 - \|\theta - \omega\|_2 > 0$ captures the additive growth in distance between parameters due to the update. To illustrate $\sigma_{\theta, \omega}^F$, we outline cases arising from aggregation rules, like SMEA, that average gradients over an adaptively chosen subset. This selection of workers can violate of the loss function's co-coercivity inequality: let S_{θ} and S_{ω} denote the subsets selected for parameters θ and ω respectively; violation occurs when $S_{\theta} \neq S_{\omega}$. The intersection and difference of the two subsets satisfy: $n - 2f \leq |S_{\theta} \cap S_{\omega}| \leq n - f$ and $|S_{\theta} \setminus S_{\omega} \cup S_{\omega} \setminus S_{\theta}| \leq 2f$. Hence, through straightforward algebraic manipulation and leveraging gradient boundedness, we obtain in the worst case scenario $\sigma_{\theta, \omega}^F \propto \frac{f}{n-f}$.

We now implement the above observation in a sketched lower bound proof. We track how the norm of the parameter difference diverges along the optimization trajectories resulting from two neighboring datasets $\mathcal{S}, \mathcal{S}'$, for T iterations of $\mathcal{A} \in \{\text{GD}, \text{SGD}\}$ with aggregation rule F and learning rate γ . We denote the parameter trajectories by $\{\theta_t\}_{t \in \{0, \dots, T-1\}}$ and $\{\theta'_t\}_{t \in \{0, \dots, T-1\}}$. We assume that the two neighboring datasets are such that the robust aggregation F does not preserve the co-coercive inequality of the loss function throughout the iterative process. For $t \in \{0, \dots, T-1\}$, define the additive terms $\sigma_t^F = \|G_{\gamma}^F(\theta_t) - G_{\gamma}^F(\theta'_t)\|_2 - \|\theta_t - \theta'_t\|_2$ and $\sigma_t^A = \frac{\gamma}{(n-f)m} \|\nabla \ell(\theta'_t; z') - \nabla \ell(\theta'_t; z)\|_2$, where z' and z are the differing samples in $\mathcal{S}', \mathcal{S}$, respectively. Then (either deterministically for GD or in expectation for SGD):

$$\|\theta_{t+1} - \theta'_{t+1}\|_2 = \|\theta_t - \theta'_t\|_2 + \sigma_t^F + \sigma_t^A.$$

We precisely exploit this construction when $F = \text{SMEA}$, designing a loss function and neighboring datasets that enable us to lower bound the additive terms as $\sigma_t^F \geq \gamma C \frac{f}{n-f}$ and $\sigma_t^A \geq \frac{\gamma C}{(n-f)m}$. This recovers the lower bound stated in the theorem, matching the convex upper bound of Theorem 3.2. The main challenge lies in designing such a scenario, which we detail in Appendices D.5 and D.6. \square

The failure of robust aggregation rules to preserve the loss function’s co-coercivity inequality underlies the additional instability observed in our analysis, compared to the simple averaging rule. Both the SMEA and CWTM rules exemplify this limitation. Notably, under the assumption that the loss function is separable across dimensions, CWTM does preserve the loss function’s co-coercive inequality (see Remark D.2), but fails to do so in more general settings (see Remark D.3). Identifying a robust aggregation rule that maintains this inequality—or proving that such preservation is impossible—remains an important open question, which we leave for future work.

3.3 Stability discrepancy between the two threat models for nonconvex loss functions

In this section, we show how the smoothness of the loss function can be leveraged to obtain sharper stability upper bounds for robust SGD under poisoning attacks—compared to the Byzantine case—in the *nonconvex setting*. This motivates CWTM (see Definition C.2), as it is an instance of an aggregation rule (see proof in Appendix D.3) that preserves the loss function’s smoothness, due to the Lipschitz continuity of the trimmed mean operation (Lemma D.6). Our focus is SGD, since GD lacks meaningful stability guarantees in the nonconvex regime issue illustrated in [23, Figure 10] and demonstrated in [11, Section 6]. This instability arises because a single differing sample affects *every* step of GD. In contrast, SGD only incorporates the differing sample when it is sampled, allowing us to reason about the first such occurrence. This key distinction enables our analysis to yield the following stability upper bound.

Theorem 3.4 (Stability under Byzantine failure attacks in the nonconvex regime - SGD). *Consider the setting described in Section 2 under Byzantine failure attacks. Let $\mathcal{A} = \text{SGD}$, with $(f, \kappa, \|\cdot\|_{op})$ -robust aggregation rule F with $\|\cdot\|_{op} \in \{\text{Tr}, \|\cdot\|_{sp}\}$. Suppose \mathcal{A} is run for $T \in \mathbb{N}^*$ iterations. Assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ bounded by ℓ_∞ , nonconvex, C -Lipschitz and L -smooth. Then, with a monotonically non-increasing learning rate $\gamma_t \leq \frac{c}{Lt}$ with $t \in \{0, \dots, T-1\}$ and $c > 0$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq 2 \left(\frac{2C^2}{L} \right)^{\frac{1}{c+1}} \left(\frac{1}{(n-f)m} + \sqrt{\kappa} \right)^{\frac{1}{c+1}} \left(\frac{\ell_\infty T}{m} \right)^{\frac{c}{c+1}}. \quad (6)$$

Theorem 3.4 recovers the known bound for $f = 0$ [28, Theorem 4.1 and discussion therein], but exhibits an added $\sqrt{\kappa}$ term when $f > 0$. As foreshadowed in the introduction of this section, when the aggregation rule preserves the loss function’s smoothness, we can derive an improved upper bound under poisoning attacks. We defer the proof of the following upper bound proof to Appendix D.4.2.

Theorem 3.5 (Stability under data poisoning attacks in the nonconvex regime - CWTM-SGD). *Consider the setting described in Section 2 under data poisoning attacks. Let $\mathcal{A} = \text{SGD}$, with CWTM. Suppose \mathcal{A} is run for $T \in \mathbb{N}^*$ iterations and that there exist $\nu > 0$ constant such that $n \geq (2 + \nu)f$ —that is, f/n is strongly bounded away from $1/2$. Assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ bounded by ℓ_∞ , nonconvex, C -Lipschitz and L -smooth. Then, with a monotonically non-increasing learning rate $\gamma_t \leq \frac{\nu}{2+\nu} \frac{c}{Lt}$ with $t \in \{0, \dots, T-1\}$ and $c > 0$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq 2 \left(\frac{2C^2 \nu^2}{(2 + \nu)^2 L} \right)^{\frac{1}{c+1}} \frac{(T \ell_\infty)^{\frac{c}{c+1}}}{m \sqrt{n}^{\frac{1}{c+1}}}. \quad (7)$$

For $f = 0$, eq. (7) yields a suboptimal \sqrt{n} dependence instead of n [28]—a proof artifact, as this gap does not appear in the Byzantine upper bound Equation (6), which also applies to the poisoning setting. For $f > 0$, let $\varepsilon^{\text{Byzantine}}$ denote the upper bound under Byzantine failures from eq. (6), and $\varepsilon_{\text{CWTM}}^{\text{poisoning}}$ the bound under poisoning attacks with CWTM from eq. (7). We compare: $\varepsilon^{\text{Byzantine}} / \varepsilon_{\text{CWTM}}^{\text{poisoning}} = \left(\frac{(2+\nu)^2}{\nu^2} \left(\frac{\sqrt{n}}{n-f} + m\sqrt{n\kappa} \right) \right)^{\frac{1}{c+1}}$. Interestingly, leveraging the smoothness-preserving property of CWTM, the bound in the data poisoning case exhibits a better dependence on m : as m grows unbounded, $\varepsilon_{\text{CWTM}}^{\text{poisoning}}$ decreases faster than $\varepsilon^{\text{Byzantine}}$. We defer additional comparisons using SMEA under data poisoning attack in Appendix D.4.3.

4 Generalization gap between byzantine Ffailures and poisoning attacks

In this section, we analyze how the stability gap influences generalization under both threat models. While [15] shows identical optimization error impact for both attacks, our stability analysis reveals a discrepancy in the generalization error. In what follows, we consider the smooth and convex setting with the SMEA aggregation rule, for which $\kappa_{\text{SMEA}} = \frac{4f}{n-f} \left(1 + \frac{f}{n-2f}\right)^2$ [5, Proposition 5.1]. First, we compare our theoretical stability bounds under Byzantine failure and data poisoning attacks, i.e.,

$$\underbrace{2\gamma C^2 T \left(2\sqrt{\frac{f}{n-f}} \left(1 + \frac{f}{n-2f} \right) + \frac{1}{(n-f)m} \right)}_{\epsilon_{\text{SMEA}}^{\text{Byzantine}} \text{ from (2)}} \quad \text{and} \quad \underbrace{2\gamma C^2 T \left(\frac{f}{n-f} + \frac{1}{(n-f)m} \right)}_{\epsilon_{\text{SMEA}}^{\text{poisoning}} \text{ from (4)}},$$

with numerical observations. Recall that the green term correspond to the uniform stability without any attack, while the red and the blue terms correspond to added instabilities due to Byzantine failure and data poisoning attacks, respectively.

Implications for the generalization error. We instantiate the problem outlined in the proof of our lower bound (Theorem 3.3, see Appendix D.5.1), under data poisoning attack, with $\gamma = C = 1$ and $m = 1$. To calculate the generalization error, we assume every honest local distribution is a Dirac with singularity at the value defined in the proof of Theorem 3.3—except for a “pivot” worker (indexed 1 without loss of generality) whose distribution is a mixture $p_1 = \frac{1}{2}\delta_0 + \frac{1}{2}\delta_{-1}$. In this case, the two possible datasets are always neighboring. We denote by $\theta_T^{(0)}$ and $\theta_T^{(-1)}$ the parameters output by the robust distributed GD algorithm with SMEA on the two different datasets. In this setup, we are able to compute the generalization error exactly as follows:

$$\mathbb{E}_{\mathcal{S}} [R_{\mathcal{H}}(\mathcal{A}(\mathcal{S})) - \widehat{R}_{\mathcal{H}}(\mathcal{A}(\mathcal{S}))] = \frac{\theta_T^{(-1)} - \theta_T^{(0)}}{4(n-f)} \text{ (see Appendix E.2 for the proof).}$$

This value corresponds to the stability plots shown in Figure 1, modulo a scaling factor of $\frac{1}{4(n-f)}$. Hence, in the worst-case, the generalization error under Byzantine failure and data poisoning attacks exhibits the same discrepancy as the one observed in uniform algorithmic stability.

Numerical validation. For our experiments, we fix the number of epochs to $T = 5$, the number of total workers $n = 15$ and vary the number of misbehaving workers f from 1 to 7. We observe that, under data poisoning, our theoretical upper bound closely matches the empirically measured algorithmic stability (see Figure 1), validating the practical tightness of our analysis in this setting. This result was expected, given our tightness proof. Next, we evaluate the system’s performance under a specific Byzantine failure attack strategy designed to maximize its influence. This attack mimics the poisoning behavior in the first iteration, then amplifies the resulting parameter direction by adaptively crafting updates that remain within the aggregation rule’s selection range (see Appendix E.1 for details). We observe that this Byzantine failure attack results in a significantly higher numerical algorithmic instability that grows with the number of Byzantine workers. The numerical observation — i.e., the solid red line in Figure 1 — aligns with our theoretical upper bound (the solid line with crosses). This agreement becomes even more apparent when using an empirical estimate of the robustness coefficient κ in place of its theoretical worst-case value (dashed line with crosses; see Appendix E.1 for details). We argue that filling this gap, while informative, is orthogonal to our main goal: assessing the applicability of our upper bound under Byzantine failure attacks. However, proving this theoretically constitutes an interesting future research direction. Overall, this finding provides strong evidence supporting the tightness of our upper bound for Byzantine failure attacks, assuming scenarios where the robustness parameter is maximized. Remarkably, this also reveals a fundamental generalization gap between Byzantine and data poisoning threat models, as highlighted earlier.

5 Related work

Minimizing empirical risk over honest workers in the presence of Byzantine workers has been extensively studied [22, 25, 26, 4, 19, 5], leading to robust aggregation schemes enabling distributed SGD to attain optimal error even under data heterogeneity. A key contribution of this line of work is

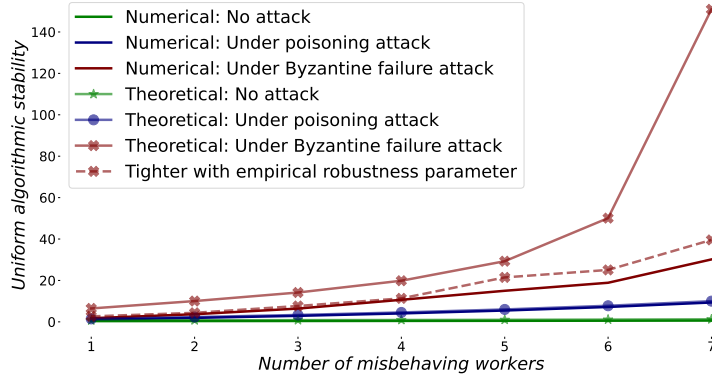


Figure 1: Stability bound for both attacks, with empirical measurements under both optimal poisoning and tailored Byzantine failure attacks. Experiments use gradient descent, so no variance is reported.

the formal analysis of aggregation rules previously studied under i.i.d. assumptions. However, these analyses focus on optimization error and offer limited insight into the generalization performance.

Empirical studies [8, 43, 38, 3, 4] have shown that the generalization error under Byzantine attacks is often worse than under data poisoning, but offer no formal explanation for this gap. These findings stand in contrast to theoretical analyses proving comparable statistical error rates across the two threat models [2, 48, 47]. However, these theoretical results assume an i.i.d. data setting, where all honest workers sample from the same distribution, and do not account for data heterogeneity common in distributed learning. We bridge this gap between empirical observations and theoretical results by developing the first theoretical framework to analyze the generalization error of robust first-order iterative methods under heterogeneous data and the two considered threat models. Our findings apply to a much larger class of loss functions, including smooth non-convex and convex, and indeed show that Byzantine attacks can be more harmful to generalization than data poisoning.

A prior work [17], which claims an equivalence between the two threat models in the context of PAC learning, considers only the classic non-robust distributed SGD algorithm. Their analysis does not apply to robust distributed GD or SGD. Another work [15], which analyzes the relationship between data poisoning and Byzantine attacks for robust distributed SGD, considers a streaming setting in which honest workers acquire i.i.d. training samples from their respective data-generating distributions at each iteration. While they show that in that setting the two threat models yield comparable population risk (asymptotically), their results and proof techniques do not apply to a more pragmatic setting that we consider wherein the training datasets across the honest workers are sampled *a priori*.

Our work relates to [44], which studies uniform stability of robust decentralized SGD under Byzantine failure attacks. However, their analysis relies on a diameter-based robustness definition for characterizing an aggregation rule. This definition leads to robustness parameters that exhibit suboptimal dependence on the problem dimension or the number of honest agents, as highlighted in [46, Table II] and further discussed in [4, Section 8.3]. Consequently, this robustness notion yields suboptimal empirical resilience (cf. [16, 4]), in contrast to covariance-based definitions, such as $(f, \kappa, \|\cdot\|_{op})$ -robustness, that we consider. Moreover, [44] does not provide any insights into the data poisoning threat model.

6 Conclusion & future work

We present the first theoretically grounded explanation for the gap in generalization performance between two central threat models in robust distributed learning: Byzantine and data poisoning attacks. Our results show that Byzantine failures can cause strictly worse generalization—even when both attacks are optimally designed. This reveals a fundamental distinction between the two threat models, overlooked by prior analyses focused solely on optimization error, and emphasizes the practical significance of our findings.

We see several promising directions for future work: (i) Identify a robust aggregation rule that preserves the loss function’s co-coercivity—or prove such preservation is impossible; (ii) Go beyond

worst-case scenarios by considering data-dependent refinements like those based on on-average stability [29, 39, 45]; (iii) Extend our results to the *locally-poisonous* setting [15]; (iv) Investigate whether the trilemma between robustness, privacy and optimization error shown in recent work [5] also holds for generalization error.

Acknowledgments and Disclosure of Funding

The work of Thomas Boudou and Aurélien Bellet is supported by grant ANR 22-PECY-0002 IPOP (Interdisciplinary Project on Privacy) project of the Cybersecurity PEPR.

References

- [1] K. Abbaszadeh, C. Pappas, J. Katz, and D. Papadopoulos. Zero-knowledge proofs of training for deep neural networks. In B. Luo, X. Liao, J. Xu, E. Kirda, and D. Lie, editors, *CCS*, 2024.
- [2] D. Alistarh, Z. Allen-Zhu, and J. Li. Byzantine stochastic gradient descent. *Advances in neural information processing systems*, 31, 2018.
- [3] Z. Allen-Zhu, F. Ebrahimi, J. Li, and D. Alistarh. Byzantine-resilient non-convex stochastic gradient descent. *arXiv preprint arXiv:2012.14368*, 2020.
- [4] Y. Allouah, S. Farhadkhani, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. Fixing by mixing: A recipe for optimal byzantine ml under heterogeneity. In F. Ruiz, J. Dy, and J.-W. van de Meent, editors, *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, volume 206 of *Proceedings of Machine Learning Research*, pages 1232–1300. PMLR, 2023.
- [5] Y. Allouah, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. On the privacy-robustness-utility trilemma in distributed learning. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 569–626. PMLR, 2023.
- [6] P. Alquier. User-friendly introduction to pac-bayes bounds. *Foundations and Trends® in Machine Learning*, 17(2):174–303, 2024.
- [7] F. Bach. *Learning Theory from First Principles*. Adaptive Computation and Machine Learning series. MIT Press, 2024.
- [8] G. Baruch, M. Baruch, and Y. Goldberg. A little is enough: Circumventing defenses for distributed learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- [9] L. Bottou and O. Bousquet. The tradeoffs of large scale learning. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007.
- [10] O. Bousquet and A. Elisseeff. Stability and generalization. *Journal of Machine Learning Research*, 2:499–526, 2002.
- [11] Z. Charles and D. Papailiopoulos. Stability and generalization of learning algorithms that converge to global optima. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 745–754. PMLR, 2018.
- [12] L. Devroye and T. Wagner. Distribution-free performance bounds for potential function rules. *Information Theory, IEEE Transactions on*, IT-25:601 – 604, 10 1979.
- [13] A. Elisseeff, T. Evgeniou, and M. Pontil. Stability of randomized learning algorithms. *Journal of Machine Learning Research*, 6(3):55–79, 2005.
- [14] S. Farhadkhani, R. Guerraoui, N. Gupta, and R. Pinot. Brief announcement: a case for byzantine machine learning. In *Proceedings of the 43rd ACM Symposium on Principles of Distributed Computing*, pages 131–134, 2024.
- [15] S. Farhadkhani, R. Guerraoui, N. Gupta, and R. Pinot. On the relevance of byzantine robust optimization against data poisoning. *arXiv preprint arXiv:2405.00491*, 2024.
- [16] S. Farhadkhani, R. Guerraoui, N. Gupta, R. Pinot, and J. Stephan. Byzantine machine learning made easy by resilient averaging of momentums. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 6246–6283. PMLR, 2022.
- [17] S. Farhadkhani, R. Guerraoui, O. Vilemaud, et al. An equivalence between data poisoning and byzantine gradient attacks. In *International Conference on Machine Learning*, pages 6284–6323. PMLR, 2022.

- [18] S. Goldwasser, S. Micali, and C. Rackoff. The knowledge complexity of interactive proof systems. *SIAM Journal on Computing*, 18(1):186–208, 1989.
- [19] E. Gorbunov, S. Horváth, P. Richtárik, and G. Gidel. Variance reduction is an antidote to byzantines: Better rates, weaker assumptions and communication compression as a cherry on the top. In *The Eleventh International Conference on Learning Representations*, 2023.
- [20] R. Guerraoui, N. Gupta, and R. Pinot. Byzantine machine learning: A primer. *ACM Comput. Surv.*, 56(7), 2024.
- [21] R. Guerraoui, N. Gupta, and R. Pinot. *Robust Machine Learning: Distributed Methods for Safe AI*. Machine Learning: Foundations, Methodologies, and Applications Series. Springer Nature Singapore, Imprint: Springer, 2024.
- [22] R. Guerraoui, S. Rouault, et al. The hidden vulnerability of distributed learning in byzantium. In *International conference on machine learning*, pages 3521–3530. PMLR, 2018.
- [23] M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International conference on machine learning*, pages 1225–1234. PMLR, 2016.
- [24] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings, R. G. L. D’Oliveira, H. Eichner, S. E. Rouayheb, D. Evans, J. Gardner, Z. Garrett, A. Gascón, B. Ghazi, P. B. Gibbons, M. Gruteser, Z. Harchaoui, C. He, L. He, Z. Huo, B. Hutchinson, J. Hsu, M. Jaggi, T. Javidi, G. Joshi, M. Khodak, J. Konečný, A. Korolova, F. Koushanfar, S. Koyejo, T. Lepoint, Y. Liu, P. Mittal, M. Mohri, R. Nock, A. Özgür, R. Pagh, H. Qi, D. Ramage, R. Raskar, M. Raykova, D. Song, W. Song, S. U. Stich, Z. Sun, A. T. Suresh, F. Tramèr, P. Vepakomma, J. Wang, L. Xiong, Z. Xu, Q. Yang, F. X. Yu, H. Yu, and S. Zhao. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021.
- [25] S. P. Karimireddy, L. He, and M. Jaggi. Learning from history for byzantine robust optimization. In *International Conference on Machine Learning*, pages 5311–5319. PMLR, 2021.
- [26] S. P. Karimireddy, L. He, and M. Jaggi. Byzantine-robust learning on heterogeneous datasets via bucketing. In *International Conference on Learning Representations*, 2022.
- [27] L. Lamport, R. Shostak, and M. Pease. The byzantine generals problem. *ACM Transactions on Programming Languages and Systems*, 4(3):382–401, 1982.
- [28] B. Le Bars, A. Bellet, M. Tommasi, K. Scaman, and G. Neglia. Improved stability and generalization guarantees of the decentralized SGD algorithm. In R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, and F. Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 26215–26240. PMLR, 2024.
- [29] Y. Lei and Y. Ying. Fine-grained analysis of stability and generalization for stochastic gradient descent. In *International Conference on Machine Learning*, pages 5809–5819. PMLR, 2020.
- [30] S. Liu, N. Gupta, and N. H. Vaidya. Approximate byzantine fault-tolerance in distributed optimization. In *Proceedings of the 2021 ACM Symposium on Principles of Distributed Computing*, pages 379–389, 2021.
- [31] G. Lugosi and G. Neu. Generalization bounds via convex analysis. In *Conference on Learning Theory*, pages 3524–3546. PMLR, 2022.
- [32] D. A. McAllester. Pac-bayesian model averaging. *Machine Learning*, 37(3):275–299, 1999.
- [33] W. Rogers and T. Wagner. A finite sample distribution-free performance bound for local discrimination rules. *The Annals of Statistics*, 6, 05 1978.
- [34] D. Russo and J. Zou. How much does your data exploration overfit? controlling bias via information usage. *IEEE Transactions on Information Theory*, 66(1):302–323, 2019.
- [35] C. Sabater, A. Bellet, and J. Ramon. An Accurate, Scalable and Verifiable Protocol for Federated Differentially Private Averaging. *Machine Learning*, 111:4249–4293, 2022.
- [36] S. Shalev-Shwartz, O. Shamir, N. Srebro, and K. Sridharan. Learnability, stability and uniform convergence. *Journal of Machine Learning Research*, 11(90):2635–2670, 2010.
- [37] A. S. Shamsabadi, G. Tan, T. I. Cebere, A. Bellet, H. Haddadi, N. Papernot, X. Wang, and A. Weller. Confidential-DPproof: Confidential Proof of Differentially Private Training. In *ICLR*, 2024.
- [38] V. Shejwalkar and A. Houmansadr. Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning. In *NDSS*, 2021.
- [39] Z. Sun, X. Niu, and E. Wei. Understanding generalization of federated learning via stability: Heterogeneity matters. In *International Conference on Artificial Intelligence and Statistics*, pages 676–684. PMLR, 2024.
- [40] J. Thaler. Proofs, arguments, and zero-knowledge. *Found. Trends Priv. Secur.*, 4(2-4):117–660, 2022.

- [41] V. Vapnik and A. Chervonenkis. *Theory of Pattern Recognition [in Russian]*. Nauka, Moscow, 1974. (German Translation: W. Wapnik & A. Tscherwonienkis, *Theorie der Zeichenerkennung*, Akademie-Verlag, Berlin, 1979).
- [42] V. N. Vapnik. *Statistical learning theory*. Wiley-interscience, 1998.
- [43] C. Xie, O. Koyejo, and I. Gupta. Fall of empires: Breaking byzantine-tolerant sgd by inner product manipulation. In *Uncertainty in Artificial Intelligence*, pages 261–270. PMLR, 2020.
- [44] H. Ye and Q. Ling. Generalization error matters in decentralized learning under byzantine attacks. *IEEE Transactions on Signal Processing*, 73:843–857, 2025.
- [45] H. Ye, T. Sun, and Q. Ling. Generalization guarantee of decentralized learning with heterogeneous data. In *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2025.
- [46] H. Ye, H. Zhu, and Q. Ling. On the tradeoff between privacy preservation and byzantine-robustness in decentralized learning. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 9336–9340, 2024.
- [47] D. Yin, Y. Chen, R. Kannan, and P. Bartlett. Byzantine-robust distributed learning: Towards optimal statistical rates. In *International conference on machine learning*, pages 5650–5659. Pmlr, 2018.
- [48] B. Zhu, L. Wang, Q. Pang, S. Wang, J. Jiao, D. Song, and M. I. Jordan. Byzantine-robust federated learning with optimal statistical rates and privacy guarantees, 2023.

Organization of the appendix

- Appendix A contains deferred proofs for Section 2, focusing on the link between uniform stability and generalization bounds under Byzantine failure and poisoning attacks.
- Appendix B summarizes the regularity assumptions and includes key expansivity results.
- Appendix C formally defines SMEA and CWTM, along with their key known properties.
- Appendix D provides detailed proofs for the stability analysis of robust distributed GD and SGD presented in Section 3. This includes: bounds on the expected (trace and spectral) norm of the empirical covariance of the honest gradients' (Appendix D.1), used for bounding uniform stability under Byzantine attacks; upper bounds on uniform stability under Byzantine attacks (Appendix D.2), with detailed proofs of Theorems 3.1 and 3.4; properties of the CWTM aggregation rule (Appendix D.3), used to improve stability bounds under poisoning attacks for smooth nonconvex objectives; upper bounds on uniform stability under poisoning attacks with SMEA and CWTM (Appendix D.4), corresponding to proofs of Theorems 3.2 and 3.5; and lower bounds on uniform stability under poisoning attacks with SMEA (Appendices D.5 and D.6), including the proof of Theorem 3.3.
- Appendix E specifies details supporting Section 4, including numerical experiments (Appendix E.1) and generalization error computations (Appendix E.2).
- Appendix F gives example of results with refined bounded heterogeneity and bounded variance assumptions instead of the more general bounded gradients assumption.
- Appendix G discusses high-probability excess risk bounds for robust distributed learning under smooth and strongly convex loss functions.

A Link between stability and generalization under Byzantine failure or poisoning attacks

In this section, we focus on the formal link between uniform stability and generalization under Byzantine and poisoning attacks. The proof follows classic arguments [10].

Lemma 2.1 (Generalization in expectation through uniform algorithmic stability with Byzantine workers). *Consider the setting described in Section 2, under either Byzantine failure or poisoning attacks. let \mathcal{A} an ε -uniformly stable (randomized) distributed algorithm. Then,*

$$|\mathbb{E}_{\mathcal{S}, \mathcal{A}}[R_{\mathcal{H}}(\mathcal{A}(\mathcal{S})) - \hat{R}_{\mathcal{H}}(\mathcal{A}(\mathcal{S}))]| \leq \varepsilon.$$

Proof. Recall that \mathcal{S} denotes the collective dataset of honest workers. Let $\mathcal{S}' = \{z'^{(i,j)}, i \in \mathcal{H}, j \in \{1, \dots, m\}\}$ be another independently sampled dataset for honest workers where for all $i \in \mathcal{H}$, $\mathcal{D}'_i = \{z'^{(i,1)}, \dots, z'^{(i,m)}\}$ is composed of m i.i.d. data points (or samples) drawn from distribution p_i . Let $\mathcal{S}^{(i,j)}$ be a neighboring dataset of \mathcal{S} with only a single differing data sample $z'^{(i,j)}$ at index $(i, j) \in \mathcal{H} \times \{1, \dots, m\}$. Then,

$$\begin{aligned} \mathbb{E}_{\mathcal{S}, \mathcal{A}}[\hat{R}_{\mathcal{H}}(\mathcal{A}(\mathcal{S}))] &=_{z^{(i,j)}, z'^{(i,j)} \text{ i.i.d.}} \mathbb{E}_{\mathcal{S}, \mathcal{S}', \mathcal{A}} \left[\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \frac{1}{m} \sum_{j=1}^m \ell(\mathcal{A}(\mathcal{S}^{(i,j)}), z'^{(i,j)}) \right] \\ &= \mathbb{E}_{\mathcal{S}, \mathcal{S}', \mathcal{A}} \left[\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \frac{1}{m} \sum_{j=1}^m \ell(\mathcal{A}(\mathcal{S}), z'^{(i,j)}) \right] + \delta_{\mathcal{A}(\mathcal{S}) \perp\!\!\!\perp z'^{(i,j)}} \mathbb{E}_{\mathcal{S}, \mathcal{A}}[R_{\mathcal{H}}(\mathcal{A}(\mathcal{S}))] + \delta \end{aligned}$$

with

$$\begin{aligned} \delta &= \mathbb{E}_{\mathcal{S}, \mathcal{S}'} \left[\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \frac{1}{m} \sum_{j=1}^m \mathbb{E}_{\mathcal{A}} \left[\ell(\mathcal{A}(\mathcal{S}^{(i,j)}), z'^{(i,j)}) - \ell(\mathcal{A}(\mathcal{S}), z'^{(i,j)}) \right] \right] \\ &\leq \mathbb{E}_{\mathcal{S}} \mathbb{E}_{\mathcal{S}'} \left[\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \frac{1}{m} \sum_{j=1}^m \sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}} \left[\ell(\mathcal{A}(\mathcal{S}^{(i,j)}), z) - \ell(\mathcal{A}(\mathcal{S}), z) \right] \right] \leq \varepsilon \end{aligned}$$

Substituting from the above equation proves the lemma. \square

B Regularity assumptions and expansivity results

To analyze the stability of iterative optimization algorithms, we typically rely on regularity assumptions on the loss function. We begin by stating standard regularity assumptions [23] used throughout the paper, along with classical results that follow from them.

Definition B.1 (Regularity assumptions). *Let $\ell : \Theta \rightarrow \mathbb{R}$ differentiable, $\Theta \subset \mathbb{R}^d$.*

(i) ℓ is C -Lipschitz-continuous (i.e., C -Lipschitz) if there exists $C > 0$ such that,

$$\forall u, v \in \Theta, \quad |\ell(u) - \ell(v)| \leq C \|u - v\|_2.$$

This property is equivalent to the norm of the gradient of ℓ being uniformly bounded by C .

(ii) ℓ is L -Lipschitz-smooth (i.e., L -smooth) if its gradient is L -Lipschitz. This is equivalent to the following smoothness inequality:

$$\forall u, v \in \Theta, \quad \ell(u) \leq \ell(v) + \langle \nabla \ell(v), u - v \rangle + \frac{L}{2} \|u - v\|_2^2. \quad (8)$$

(iii) ℓ is convex if and only if, $\forall u, v \in \Theta$, $\langle \nabla \ell(u) - \nabla \ell(v), u - v \rangle \geq 0$.

Moreover it is μ -strongly convex if there exists $\mu > 0$ such that,

$$\forall u, v \in \Theta, \quad \ell(u) \geq \ell(v) + \langle \nabla \ell(v), u - v \rangle + \frac{\mu}{2} \|u - v\|_2^2.$$

We note that for a strongly convex function to have bounded gradients, it must be defined on a convex compact set, which then implies boundedness of both the loss and the gradient. To ensure this, we can either penalize the problem or restrict the parameter domain and apply an Euclidean projection at every step. Throughout the paper, we will tacitly assume this when referring to strongly convex functions, which does not limit the applicability of our analysis since the projection does not increase the distance between projected points.

Notably, convex and L -smooth functions satisfy an important property known as co-coercivity (see Proposition 5.4 in [7]). Specifically,

Lemma B.1 (Co-coercivity). *Let $\ell : \Theta \rightarrow \mathbb{R}$ differentiable, $\Theta \subset \mathbb{R}^d$. If ℓ is a convex and L -smooth function, then it satisfies the co-coercivity inequality,*

$$\forall u, v \in \Theta, \quad \langle \nabla \ell(u) - \nabla \ell(v), u - v \rangle \geq \frac{1}{L} \|\nabla \ell(u) - \nabla \ell(v)\|_2^2. \quad (9)$$

We focus on robust variants of distributed gradient descent and stochastic gradient descent. Below we recall properties of the one-step update. Let $\mathcal{A} \in \{\text{GD}, \text{SGD}\}$, $\ell : \Theta \rightarrow \mathbb{R}$ differentiable and $\gamma > 0$, we denote:

$$G_\gamma^{\mathcal{A}}(\theta) := \theta - \gamma \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} g^{(i)}.$$

Here, for each $i \in \mathcal{H}$, $g^{(i)}$ denotes the exact gradient $\nabla \ell(\theta)$ in the case of GD, or an unbiased estimate thereof—effectively computed with a sample from the local dataset—in the case of SGD. We recall the expansiveness definition for an update rule from [23, Definition 2.3]:

Definition B.2 (η -expansive update rule). *An update rule $G : \Theta \rightarrow \Theta$ is said to be η -expansive if for all $\theta, \omega \in \Theta$,*

$$\|G(\theta) - G(\omega)\|_2 \leq \eta \|\theta - \omega\|_2.$$

A straightforward extension of [23, Lemma 3.6] to the distributed setting yields:

Lemma B.2 (Expansivity of gradient updates). *Let $\mathcal{A} \in \{\text{GD}, \text{SGD}\}$ and $\gamma > 0$. Assume $\forall z \in \mathcal{Z}, \ell(\cdot, z) : \Theta \rightarrow \mathbb{R}$ is L -smooth, Then $G_\gamma^{\mathcal{A}}$ is $\eta_{G_\gamma^{\mathcal{A}}}$ -expansive [23, Definition 2.3], with the following expressions:*

(i) $\eta_{G_\gamma^{\mathcal{A}}} = (1 + \gamma L)$.

(ii) Assume in addition that ℓ is convex. Then, for any $\gamma \leq \frac{2}{L}$, $\eta_{G_\gamma^{\mathcal{A}}} = 1$.

- (iii) Assume in addition that ℓ is μ -strongly convex. Then, for any $\gamma \leq \frac{2}{\mu+L}$, $\eta_{G_\gamma^A} = (1 - \gamma \frac{L\mu}{\mu+L})$. For $\gamma \leq \frac{1}{L}$, since we have $L \geq \mu$, we can simplify the contractive constant to $\eta_{G_\gamma^A} = (1 - \gamma\mu)$.

$\eta_{G_\gamma^A}$ is referred to as the expansivity coefficient of \mathcal{A} with learning rate γ .

C Robust aggregation rules

We present the *Component-Wise Trimmed Mean* (CWTM) and the *Smallest Maximum Eigenvalue Averaging* (SMEA) that we use throughout the paper. They are, respectively, instances of (f, κ, Tr) -robust and $(f, \kappa, \|\cdot\|_{sp})$ -robust aggregation rules, each with an asymptotically optimal robustness parameter. Note that (f, κ, Tr) -robustness is computationally cheaper to satisfy and yields optimal convergence error in settings where the honest gradients' covariance matrix has uniformly bounded trace over the parameter space. On the other hand, $(f, \kappa, \|\cdot\|_{sp})$ -robustness is computationally expensive to satisfy—as robust aggregation rules must effectively operate in high-dimensional spaces, unlike simpler methods such as CWTM—but yields optimal convergence error when only the spectral norm of the honest gradients' covariance matrix is assumed to be bounded.

Definition C.1 (Smallest Maximum Eigenvalue Averaging (SMEA)). *Given $f, n \in \mathbb{N}$, $f < n/2$ and vectors $g_1, \dots, g_n \in \mathbb{R}^d$, the smallest maximum eigenvalue averaging aggregation rule outputs:*

$$\bar{g}_{S^*} = \frac{1}{|S^*|} \sum_{i \in S^*} g_i$$

$$\text{with } S^* \in \underset{S \subseteq \{1, \dots, n\}, |S|=n-f}{\text{argmin}} \lambda_{\max} \left(\frac{1}{|S|} \sum_{i \in S} (g_i - \bar{g}_S)(g_i - \bar{g}_S)^\top \right)$$

The SMEA aggregation rule is especially valuable because it is $(f, \kappa, \|\cdot\|_{sp})$ -robust, with an optimal asymptotic robustness coefficient.

Lemma C.1 (Robustness of SMEA [5, Prop. 5.1]). *SMEA is $(f, \kappa, \|\cdot\|_{sp})$ -robust with*

$$\kappa = \frac{4f}{n-f} \left(1 + \frac{f}{n-2f} \right)^2.$$

Assume that there exist ν constant such that $n \geq (2 + \eta)f$, that is f/n is strongly away from $1/2$, then $\kappa \in \mathcal{O}(f/n)$ is asymptotically optimal.

Unlike SMEA, which leverages spectral properties to achieve robust aggregation, CWTM relies on coordinate-wise trimming.

Definition C.2 (Component-Wise Trimmed Mean (CWTM)). *Given $f, n \in \mathbb{N}$, $f < n/2$, and real values x_1, \dots, x_n , we denote by σ a permutation on $\{1, \dots, n\}$ that sorts the values in non-decreasing order $x_{\sigma(1)} \leq \dots \leq x_{\sigma(n)}$. Let $S_x = \{i \in [n]; f+1 \leq \sigma(i) \leq n-f\}$, the trimmed mean is given by:*

$$TM(x_1, \dots, x_n) = \frac{1}{n-2f} \sum_{i=f+1}^{n-f} x_{\sigma(i)} = \frac{1}{n-2f} \sum_{i \in S_x} x_i$$

We define the component-wise trimmed mean for vectors $x_1, \dots, x_n \in \mathbb{R}^d$ as the coordinate-wise application of the trimmed mean operation. That is $\text{CWTM}(x_1, \dots, x_n)$ is the vector whose k -th coordinate is defined as

$$[\text{CWTM}(x_1, \dots, x_n)]_k = TM([x_1]_k, \dots, [x_n]_k)$$

The CWTM aggregation rule is particularly useful as it is (f, κ, Tr) -robust, with an optimal asymptotic robustness coefficient.

Lemma C.2 (Robustness of CWTM [4, Prop. 2.1]). *CWTM is (f, κ, Tr) -robust with*

$$\kappa = \frac{6f}{n-2f} \left(1 + \frac{f}{n-2f} \right).$$

Assume that there exist ν constant such that $n \geq (2 + \eta)f$, that is f/n is strongly away from $1/2$, then $\kappa \in \mathcal{O}(f/n)$ is asymptotically optimal.

D Deferred proofs from Section 3

This appendix contains the detailed proofs and technical results supporting the analysis presented in Section 3. We begin with supporting results in Appendix D.1, specifically, we bound the expectation of both the trace and the spectral norm of the empirical covariance matrix of honest workers' gradients. Building on this, we establish bounds on uniform stability under Byzantine attacks in Appendix D.2, with detailed proofs of Theorems 3.1 and 3.4. These bounds are established for robust distributed GD and SGD using any $(f, \kappa, \|\cdot\|_{op})$ -robust aggregation rule, across strongly convex, convex, and nonconvex optimization settings. In Appendix D.3, we detail the properties of the CWTM aggregation rule, which is used to improve stability bound under poisoning attacks for smooth nonconvex objectives. Appendix D.4 focuses on leveraging loss function regularities under data poisoning, providing stability bounds for robust distributed GD and SGD using the SMEA and CWTM rules across strongly convex, convex, and nonconvex settings. This corresponds to proofs of Theorems 3.2 and 3.5. Appendices D.5 and D.6 details the unavoidable instability introduced by the SMEA aggregation rule by providing lower bounds on uniform stability under poisoning attacks for robust distributed GD and SGD across strongly convex and convex settings, that is proof of Theorem 3.3.

D.1 Expectation bound for the honest gradients' empirical covariance trace and spectral norm

We derive bounds on the expected trace and spectral norm of the empirical covariance matrix of the honest workers' gradients under two settings the bounded gradient assumption (*i.e.* Lipschitz-continuity of the loss function).

Lemma D.1 (Expectation bound for the empirical covariance trace). *Consider the setting described in Section 2 under Byzantine attacks. Assume the loss function C -Lipschitz. Let $\mathcal{A} \in \{\text{GD}, \text{SGD}\}$, with a (f, κ, Tr) -robust aggregation rule, for $T \in \mathbb{N}^*$ iterations. We have for every $t \in \{1, \dots, T-1\}$:*

$$\mathbb{E}_{\mathcal{A}} \text{Tr} \Sigma_{\mathcal{H},t} = \mathbb{E}_{\mathcal{A}} \left[\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \|g_t^{(i)} - \bar{g}_t\|_2^2 \right] \leq C^2.$$

Proof. For the first inequality, we directly use the boundedness of gradients:

$$\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \|g_t^{(i)} - \bar{g}_t\|_2^2 \leq \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \|g_t^{(i)}\|_2^2 \leq C^2.$$

□

Lemma D.2 (Expectation bound for the empirical covariance spectral norm). *Consider the setting described in Section 2 under Byzantine attacks. Assume the loss function C -Lipschitz. Let $\mathcal{A} \in \{\text{GD}, \text{SGD}\}$, with a $(f, \kappa, \|\cdot\|_{sp})$ -robust aggregation rule, for $T \in \mathbb{N}^*$ iterations. We have for every $t \in \{1, \dots, T-1\}$:*

$$\mathbb{E}_{\mathcal{A}} \|\Sigma_{\mathcal{H},t}\|_{sp} = \mathbb{E}_{\mathcal{A}} \left[\lambda_{\max} \left(\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} (g_t^{(i)} - \bar{g}_t)(g_t^{(i)} - \bar{g}_t)^\top \right) \right] \leq C^2.$$

Proof. For the first inequality, we directly use the boundedness of gradients:

$$\begin{aligned} \lambda_{\max} \left(\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} (g_t^{(i)} - \bar{g}_t)(g_t^{(i)} - \bar{g}_t)^\top \right) &= \sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, g_t^{(i)} - \bar{g}_t \rangle^2 \\ &\leq \sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, g_t^{(i)} \rangle^2 \leq \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \|g_t^{(i)}\|_2^2 \leq C^2. \end{aligned}$$

□

Remark D.1. In Appendix F, we present example results derived under refined assumptions—namely, bounded heterogeneity and bounded variance—in place of the more general bounded gradients condition (*i.e.* Lipschitz-continuity of the loss function). These assumptions are more commonly used in the optimization literature (as opposed to the generalization literature), and while they often yield tighter bounds, they do not provide additional conceptual insights within the scope of our discussion.

D.2 Proofs of Theorems 3.1 and 3.4

In this section, we present proofs for Theorems 3.1 and 3.4. Specifically, we derive upper bounds on the uniform stability of learning algorithms under Byzantine attacks, considering convex, strongly convex, and nonconvex objectives. We begin with a unified *proof sketch* that outlines the key ideas common to both theorems. The proof of Theorem 3.1 is divided into two parts: one for $\mathcal{A} = \text{GD}$, provided in Theorem D.3, and one for $\mathcal{A} = \text{SGD}$, detailed in points (i) and (ii) of Theorem D.5. While the results are presented in a unified manner in the main text—due to their similar structure and identical bounds—we separate the proofs into two distinct components for clarity. The proof of Theorem 3.4 is established in point (iii) of Theorem D.5.

Proof sketch. Let denote $\{\theta_t\}_{t \in \{0, \dots, T-1\}}$ and $\{\theta'_t\}_{t \in \{0, \dots, T-1\}}$ the optimization trajectories resulting from two neighboring datasets $\mathcal{S}, \mathcal{S}'$, for T iterations of $\mathcal{A} \in \{\text{GD}, \text{SGD}\}$ with F a $(f, \kappa, \|\cdot\|_{op})$ -robust aggregation rule, and learning rate schedule $\{\gamma_t\}_{t \in \{0, \dots, T-1\}}$. Let denote $m_t^{(i)}$ and $m_t'^{(i)} \in \mathbb{R}^d$ the local updates provided by worker $i \in \{1, \dots, n\}$ along the trajectories $t \in \{0, \dots, T-1\}$. In the context of the uniform algorithmic stability analysis via parameter sensitivity, we track how the parameters diverge along these optimization trajectories. To do so, we decompose the analysis by introducing an intermediate comparison between the robust update $G_{\gamma_t}^F$ and the averaging over honest workers update $G_{\gamma_t}^A$. This enables us to leverage either the regularity of the loss function or the robustness property of the aggregation rule. This approach is primarily motivated by the fact that Byzantine vectors cannot be assumed to exhibit regularities when comparing them.

$$\begin{aligned} \mathbb{E}_{\mathcal{A}} [\|G_{\gamma_t}^F(\theta_t) - G_{\gamma_t}^{F'}(\theta'_t)\|_2] &= \mathbb{E}_{\mathcal{A}} [\|G_{\gamma_t}^F(\theta_t) - G_{\gamma_t}^A(\theta_t) + G_{\gamma_t}^A(\theta_t) - G_{\gamma_t}^{A'}(\theta'_t) \\ &\quad + G_{\gamma_t}^{A'}(\theta'_t) - G_{\gamma_t}^{F'}(\theta'_t)\|_2] \\ &\leq \underbrace{\mathbb{E}_{\mathcal{A}} \|G_{\gamma_t}^A(\theta_t) - G_{\gamma_t}^{A'}(\theta'_t)\|_2}_A \\ &\quad + \underbrace{\gamma_t \mathbb{E}_{\mathcal{A}} [\|G_{\gamma_t}^F(\theta_t) - G_{\gamma_t}^A(\theta_t)\|_2 + \|G_{\gamma_t}^{F'}(\theta'_t) - G_{\gamma_t}^{A'}(\theta'_t)\|_2]}_B \end{aligned}$$

We then bound the term A with classical stability tools from [23], and the term B with the $(f, \kappa, \|\cdot\|_{op})$ -robust property, the Jensen inequality and the Lemma D.1 (if $\|\cdot\|_{op} = \text{Tr}$) or Lemma D.2 (if $\|\cdot\|_{op} = \|\cdot\|_{sp}$):

$$B \leq 2\gamma_t \mathbb{E}_{\mathcal{A}} \sqrt{\kappa \|\Sigma_{\mathcal{H}, t}\|_{op}} \leq 2\gamma_t \sqrt{\kappa \mathbb{E}_{\mathcal{A}} \|\Sigma_{\mathcal{H}, t}\|_{op}} \leq 2\gamma_t \sqrt{\kappa} C.$$

D.2.1 Proof of Theorem 3.1 for GD

Theorem D.3 (Stability under Byzantine failure attacks - GD). *Consider the setting described in Section 2 under Byzantine attacks. Let $\mathcal{A} = \text{GD}$, with F a $(f, \kappa, \|\cdot\|_{op})$ -robust aggregation rule, with $\|\cdot\|_{op} \in \{\text{Tr}, \|\cdot\|_{sp}\}$, for $T \in \mathbb{N}^*$ iterations. Suppose $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ C -Lipschitz and L -smooth.*

- (i) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ convex, with constant learning rate $\gamma \leq \frac{2}{L}$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}} [\|\ell(\mathcal{A}(\mathcal{S}'); z) - \ell(\mathcal{A}(\mathcal{S}); z)\|] \leq 2\gamma C^2 T \left(\frac{1}{(n-f)m} + \sqrt{\kappa} \right).$$

- (ii) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ μ -strongly convex, with constant learning rate $\gamma \leq \frac{1}{L}$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}} [\|\ell(\mathcal{A}(\mathcal{S}'); z) - \ell(\mathcal{A}(\mathcal{S}); z)\|] \leq \frac{2C^2}{\mu} \left(\frac{1}{(n-f)m} + \sqrt{\kappa} \right).$$

Proof. For each $t \in \{0, \dots, T-1\}$, following the proof sketch in D.2, we bound the term $\|G_{\gamma_t}^{\text{GD}}(\theta_t) - G_{\gamma_t}^{\text{GD}'}(\theta'_t)\|_2$ where $G_{\gamma_t}^{\text{GD}}(\theta_t) = \theta_t - \gamma_t \nabla \widehat{R}_{\mathcal{H}}(\theta_t)$. For $t \in \{0, \dots, T-1\}$, we have,

$$\|G_{\gamma_t}^{\text{GD}}(\theta_t) - G_{\gamma_t}^{\text{GD}'}(\theta'_t)\|_2 = \|G_{\gamma_t}^{\text{GD}}(\theta_t) - G_{\gamma_t}^{\text{GD}}(\theta'_t)\|_2$$

$$\begin{aligned}
& + \frac{\gamma_t}{|\mathcal{H}|m} \|\nabla \ell(\theta'_t; z'^{(a,b)}) - \nabla \ell(\theta'_t; z^{(a,b)})\|_2 \\
& \leq \eta_{G_{\gamma_t}^{\text{GD}}} \|\theta_t - \theta'_t\|_2 + \frac{2\gamma_t C}{(n-f)m}
\end{aligned}$$

where $\eta_{G_{\gamma_t}^{\text{GD}}}$ denotes the expansive coefficient, which depends on the regularity properties of the loss function (see Lemma B.2). Let define $\sigma_{\text{GD}}^F = 2C \left(\frac{1}{(n-f)m} + \sqrt{\kappa} \right)$, we then resolve the recursion by incorporating the impact of the Byzantine failure attack as demonstrated in the previous proof sketch.

(i) If the loss function is Lipschitz, smooth, and convex, then for all $t \in \{0, \dots, T-1\}$:

$$\|\theta_{t+1} - \theta'_{t+1}\|_2 \leq \|\theta_t - \theta'_t\|_2 + \gamma \sigma_{\text{GD}}^F$$

To conclude the proof, we assume a constant learning rate $\gamma \leq \frac{1}{L}$ and sum the inequality over $t \in 0, \dots, T-1$, noting that $\theta_0 = \theta'_0$. We then invoke the Lipschitz continuity assumption.

(ii) If the loss function is Lipschitz, smooth, and strongly convex, then for all $t \in \{0, \dots, T-1\}$:

$$\|\theta_{t+1} - \theta'_{t+1}\|_2 \leq (1 - \gamma\mu) \|\theta_t - \theta'_t\|_2 + \gamma \sigma_{\text{GD}}^F$$

Considering $\gamma \leq \frac{1}{L}$ and summing over $t \in \{0, \dots, T-1\}$ with $\theta_0 = \theta'_0$:

$$\|\theta_T - \theta'_T\|_2 \leq \gamma \sigma_{\text{GD}}^F \sum_{t=0}^{T-1} (1 - \gamma\mu)^t = \frac{\sigma_{\text{GD}}^F}{\mu} \left[1 - (1 - \gamma\mu)^T \right] \leq \frac{\sigma_{\text{GD}}^F}{\mu}$$

We conclude the derivation of the bound by invoking the Lipschitz continuity assumption. \square

D.2.2 Proof of Theorem 3.1 for SGD and Theorem 3.4

Unlike GD—where the differing samples affect every iteration—in SGD, it may not be sampled immediately. This distinction enables our analysis to yield better stability upper bound with SGD. Below, we state a lemma reasoning about the first occurrence of the differing samples.

Lemma D.4. *Consider the setting described in Section 2, under either Byzantine failure or poisoning attacks. We assume for all $z \in \mathcal{Z}$, $\ell(\cdot; z)$ nonnegative and uniformly bounded by $\ell_\infty \in \mathbb{R}_+$. Let \mathcal{S} and \mathcal{S}' being two neighboring dataset, we denote $\theta_T = \mathcal{A}(\mathcal{S})$, $\theta'_T = \mathcal{A}(\mathcal{S}')$ and intermediate state of the algorithm θ_t , θ'_t and $\delta_t = \|\theta_t - \theta'_t\|_2$ for any $t \in \{0, \dots, T\}$. Then for every $z \in \mathcal{Z}$ and every $t_0 \in \{0, \dots, \min(m, T)\}$, the following holds for SGD with sampling with replacement:*

$$\mathbb{E}_{\mathcal{A}}[\|\ell(\theta_T; z) - \ell(\theta'_T; z)\|] \leq \min\{\ell_\infty, \frac{t_0}{m} \ell_\infty + \mathbb{E}_{\mathcal{A}}[\|\ell(\theta_T; z) - \ell(\theta'_T; z)\| \delta_{t_0} = 0]\}$$

Proof. Let $z \in \mathcal{Z}$ and $t_0 \in \{0, \dots, \min(m, T)\}$, we have:

$$\begin{aligned}
\mathbb{E}_{\mathcal{A}}[\|\ell(\theta_T; z) - \ell(\theta'_T; z)\|] &= \mathbb{E}_{\mathcal{A}}[\|\ell(\theta_T; z) - \ell(\theta'_T; z)\| \delta_{t_0} = 0] \mathbb{P}(\delta_{t_0} = 0) \\
&\quad + \mathbb{E}_{\mathcal{A}}[\|\ell(\theta_T; z) - \ell(\theta'_T; z)\| \delta_{t_0} \neq 0] \mathbb{P}(\delta_{t_0} \neq 0) \\
&\leq \mathbb{E}_{\mathcal{A}}[\|\ell(\theta_T; z) - \ell(\theta'_T; z)\| \delta_{t_0} = 0] + \mathbb{P}(\delta_{t_0} \neq 0) \ell_\infty
\end{aligned}$$

To bound $\mathbb{P}(\delta_{t_0} \neq 0)$, we consider T_0 , the random variable of the first step the differing samples are drawn. We have:

$$\mathbb{P}(\delta_{t_0} \neq 0) \leq \mathbb{P}(T_0 \leq t_0) \leq \min(1, \frac{t_0}{m})$$

For the first inequality, we used the fact that under identical conditions (i.e., the same seed and initial state), the algorithm operating on two neighboring datasets will produce the same intermediate states until the differing samples are drawn. For the second inequality, we use the union bound to obtain $\mathbb{P}(T_0 \leq t_0) \leq \sum_{t=1}^{t_0} \mathbb{P}(T_0 = t) = \min(1, \frac{t_0}{m})$. \square

Building on the above lemma, we now state bounds for SGD.

Theorem D.5 (Stability under Byzantine failure attacks - SGD). *Consider the setting described in Section 2 under Byzantine attacks. Let $\mathcal{A} = \text{SGD}$, with F a $(f, \kappa, \|\cdot\|_{op})$ -robust aggregation rule, with $\|\cdot\|_{op} \in \{\text{Tr}, \|\cdot\|_{sp}\}$, for $T \in \mathbb{N}^*$ iterations. Suppose $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ C -Lipschitz and L -smooth.*

- (i) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ convex, with constant learning rate $\gamma \leq \frac{2}{L}$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)] \leq 2\gamma C^2 T \left(\frac{1}{(n-f)m} + \sqrt{\kappa} \right).$$

- (ii) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ μ -strongly convex, with constant learning rate $\gamma \leq \frac{1}{L}$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[\ell(\mathcal{A}(\mathcal{S}'); z) - \ell(\mathcal{A}(\mathcal{S}'); z)] \leq \frac{2C^2}{\mu} \left(\frac{1}{(n-f)m} + \sqrt{\kappa} \right).$$

- (iii) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ nonconvex, bounded by $\ell_{\infty} = \sup_{\theta \in \Theta, z \in \mathcal{Z}} \ell(\theta; z)$, with a monotonically non-increasing learning rate $\gamma_t \leq \frac{c}{Lt}$, $t \in \{0, \dots, T-1\}$, $c > 0$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)] \leq 2 \left(\frac{2C^2}{L} \right)^{\frac{1}{c+1}} \left(\frac{1}{(n-f)m} + \sqrt{\kappa} \right)^{\frac{1}{c+1}} \left(\frac{\ell_{\infty} T}{m} \right)^{\frac{c}{c+1}}.$$

Proof. For each $t \in \{1, \dots, T\}$, following the proof sketch in Appendix D.2, we bound the term $\mathbb{E}_{\mathcal{A}}[G_{\gamma_t}^{\text{SGD}}(\theta_t) - G_{\gamma_t}^{\text{SGD}'}(\theta'_t)]_2$ where $G_{\gamma_t}^{\text{SGD}}(\theta_t) = \theta_t - \gamma_t \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} g_t^{(i)}$, and for each $i \in \mathcal{H}$, $g_t^{(i)} = \nabla \ell(\theta_t; z_t^{(i, J_t^{(i)})})$ with $J_t^{(i)}$ sampled uniformly from $\{1, \dots, m\}$. We establish the bound by conditioning on whether the differing samples are selected at a given step, where $a, b \in \mathcal{H} \times \{1, \dots, m\}$ denote the indices of the differing samples. At each iteration, a single sample is selected uniformly at random from each honest worker. Consequently, with probability $(1 - \frac{1}{m})$, the differing samples are not selected:

$$\|G_{\gamma_t}^{\text{SGD}}(\theta_t) - G_{\gamma_t}^{\text{SGD}'}(\theta'_t)\|_2 = \|G_{\gamma_t}^{\text{SGD}}(\theta_t) - G_{\gamma_t}^{\text{SGD}}(\theta'_t)\|_2 \leq \eta_{G_{\gamma_t}^{\text{SGD}}} \|\theta_t - \theta'_t\|_2$$

where $\eta_{G_{\gamma_t}^{\text{SGD}}}$ denotes the expansive coefficient, which depends on the regularity properties of the loss function (see Lemma B.2). And with probability $\frac{1}{m}$, the differing samples are selected:

$$\begin{aligned} \|G_{\gamma_t}^{\text{SGD}}(\theta_t) - G_{\gamma_t}^{\text{SGD}'}(\theta'_t)\|_2 &= \|G_{\gamma_t}^{\text{SGD}}(\theta_t) - G_{\gamma_t}^{\text{SGD}}(\theta'_t)\|_2 \\ &\quad + \frac{\gamma_t}{|\mathcal{H}|} \|\nabla \ell(\theta'_t; z'^{(a,b)}) - \nabla \ell(\theta'_t; z^{(a,b)})\|_2 \\ &\leq \eta_{G_{\gamma_t}^{\text{SGD}}} \|\theta_t - \theta'_t\|_2 + \frac{2\gamma_t C}{|\mathcal{H}|}. \end{aligned}$$

We thus obtain the bound:

$$\mathbb{E}_{\mathcal{A}}[\|G_{\gamma_t}^{\text{SGD}}(\theta_t) - G_{\gamma_t}^{\text{SGD}'}(\theta'_t)\|_2] \leq \eta_{G_{\gamma_t}^{\text{SGD}}} \mathbb{E}_{\mathcal{A}}[\|\theta_t - \theta'_t\|_2] + \frac{2\gamma_t C}{(n-f)m} \quad (10)$$

Let define $\sigma_{\text{SGD}}^F = 2C \left(\frac{1}{(n-f)m} + \sqrt{\kappa} \right)$, we then resolve the recursion by incorporating the impact of the Byzantine failure attack as demonstrated in the previous proof sketch.

- (i) If the loss function is Lipschitz, smooth, and convex, then for all $t \in \{0, \dots, T-1\}$:

$$\mathbb{E}_{\mathcal{A}}[\|\theta_{t+1} - \theta'_{t+1}\|_2] \leq \mathbb{E}_{\mathcal{A}}[\|\theta_t - \theta'_t\|_2] + \gamma_t \sigma_{\text{SGD}}^F$$

To conclude the proof, we assume a constant learning rate $\gamma \leq \frac{1}{L}$ and sum the inequality over $t \in 0, \dots, T-1$, noting that $\theta_0 = \theta'_0$. We then invoke the Lipschitz continuity assumption.

(ii) If the loss function is Lipschitz, smooth, and strongly convex, then for all $t \in \{0, \dots, T-1\}$:

$$\mathbb{E}_{\mathcal{A}} \|\theta_{t+1} - \theta'_{t+1}\|_2 \leq (1 - \gamma\mu) \mathbb{E}_{\mathcal{A}} \|\theta_t - \theta'_t\| + \gamma\sigma_{\text{SGD}}^F$$

Considering $\gamma \leq \frac{1}{L}$ and summing over $t \in \{0, \dots, T-1\}$ with $\theta_0 = \theta'_0$:

$$\mathbb{E}_{\mathcal{A}} \|\theta_T - \theta'_T\|_2 \leq \gamma\sigma_{\text{SGD}}^F \sum_{t=0}^{T-1} (1 - \gamma\mu)^t = \frac{\sigma_{\text{SGD}}^F}{\mu} \left[1 - (1 - \gamma\mu)^T \right] \leq \frac{\sigma_{\text{SGD}}^F}{\mu}$$

We conclude the derivation of the bound by invoking the Lipschitz continuity assumption.

(iii) If the loss function is Lipschitz, smooth, and nonconvex, then considering $t_0 \in \{0, \dots, T-1\}$, for all $t \in \{t_0, \dots, T-1\}$:

$$\begin{aligned} \mathbb{E}_{\mathcal{A}} [\delta_{t+1} | \delta_{t_0} = 0] &\leq (1 + \gamma_t L) \mathbb{E}_{\mathcal{A}} [\delta_t | \delta_{t_0} = 0] + \gamma_t \sigma_{\text{SGD}}^F \\ &\leq_{1+x \leq e^x} e^{\gamma_t L} \mathbb{E}_{\mathcal{A}} [\delta_t | \delta_{t_0} = 0] + \gamma_t \sigma_{\text{SGD}}^F \end{aligned}$$

Let $c > 0$ and consider a monotonically non-increasing learning rate $\gamma_t \leq \frac{c}{Lt}$. We proceed as follows:

$$\mathbb{E}_{\mathcal{A}} [\delta_T | \delta_{t_0} = 0] \leq \frac{c\sigma_{\text{SGD}}^F}{L} \sum_{t=t_0}^{T-1} \frac{1}{t} \prod_{s=t+1}^{T-1} e^{\frac{c}{s}},$$

where:

$$\sum_{t=t_0}^{T-1} \frac{1}{t} \prod_{s=t+1}^{T-1} e^{\frac{c}{s}} = \sum_{t=t_0}^{T-1} \frac{1}{t} e^{c \sum_{s=t+1}^{T-1} \frac{1}{s}} \leq \sum_{t=t_0}^{T-1} \frac{1}{t} e^{c \log(\frac{T}{t})} = T^c \sum_{t=t_0}^{T-1} t^{-c-1} \leq \frac{1}{c} \left(\frac{T}{t_0} \right)^c.$$

Where we used the following sum-integral inequality: for any non-increasing and integrable function f , and for $a, b \in \mathbb{N}$, it holds that $\sum_{t=a}^b f(t) \leq \sum_{t=a}^b \int_{t-1}^t f(s) ds = \int_{a-1}^b f(s) ds$. Hence with $\alpha \in \mathbb{R}_+^*$:

$$\sum_{s=t+1}^{T-1} \frac{1}{s} \leq \int_t^{T-1} \frac{1}{s} ds = \ln\left(\frac{T-1}{t}\right) \leq \ln\left(\frac{T}{t}\right),$$

and

$$\sum_{t=t_0}^{T-1} t^{-\alpha-1} \leq \int_{t_0}^T t^{-\alpha-1} dt \leq -\frac{1}{\alpha} (T^{-\alpha} - t_0^{-\alpha}) \leq \frac{1}{\alpha} t_0^{-\alpha}.$$

Substituting our derivation into the bound from Lemma D.4 yields:

$$\mathbb{E}_{\mathcal{A}} [|\ell(\theta_T; z) - \ell(\theta'_T; z)|] \leq \frac{t_0 \ell_{\infty}}{m} + \frac{\sigma_{\text{SGD}}^F C}{L} \left(\frac{T}{t_0} \right)^c$$

We approximately minimize the expression with respect to t_0 by balancing the two terms, *i.e.*, setting them equal:

$$0 \leq \tilde{t}_0 = \left(\frac{m\sigma_{\text{SGD}}^F C}{L\ell_{\infty}} \right)^{\frac{1}{c+1}} T^{\frac{c}{c+1}} \leq \begin{matrix} \text{for sufficiently large } T \geq \frac{m\sigma_{\text{SGD}}^F C}{L\ell_{\infty}} \end{matrix} T.$$

Finally, we obtain the result by directly substituting \tilde{t}_0 into the bound:

$$\mathbb{E}_{\mathcal{A}} [|\ell(\theta_T; z) - \ell(\theta'_T; z)|] \leq 2 \left(\frac{\sigma_{\text{SGD}}^F C}{L} \right)^{\frac{1}{c+1}} \left(\frac{\ell_{\infty} T}{m} \right)^{\frac{c}{c+1}}. \quad \square$$

D.3 CWTM aggregation rule properties

We begin by analyzing the Lipschitz continuity of the Trimmed Mean (TM) robust aggregation on the real line, which is useful for showing the expansivity property of the Component-Wise Trimmed Mean (CWTM) robust update.

Lemma D.6. *The trimmed mean operation is $\frac{1}{n-2f}$ -Lipschitz continuous with respect to the $\|\cdot\|_1$ norm, and $\frac{\sqrt{n}}{n-2f}$ -Lipschitz continuous with respect to the $\|\cdot\|_2$ norm.*

Proof. We recall the definition of the trimmed mean operation: given $f, n \in \mathbb{N}$, $f < n/2$, and real values x_1, \dots, x_n , we denote by σ a permutation on $\{1, \dots, n\}$ that sorts the values in non-decreasing order $x_{\sigma(1)} \leq \dots \leq x_{\sigma(n)}$. Let $S_x = \{i \in [n]; f+1 \leq \sigma(i) \leq n-f\}$, the trimmed mean is given by:

$$TM(x_1, \dots, x_n) = \frac{1}{n-2f} \sum_{i=f+1}^{n-f} x_{\sigma(i)} = \frac{1}{n-2f} \sum_{i \in S_x} x_i$$

We first consider two sets of real numbers that differ in exactly one entry $\iota \in \{1, \dots, n\}$: $x_1, \dots, x_\iota, \dots, x_n$ and $x'_1, \dots, x'_\iota, \dots, x'_n$ such that $x_j = x'_j$ for all $j \neq \iota$. Let S_x and $S_{x'}$ denote the sets of selected indices in each case. We necessarily have $|S_{x'} \setminus S_x| \leq 1$. In the figure below, we illustrate different scenarios to provide an intuition for our proof.

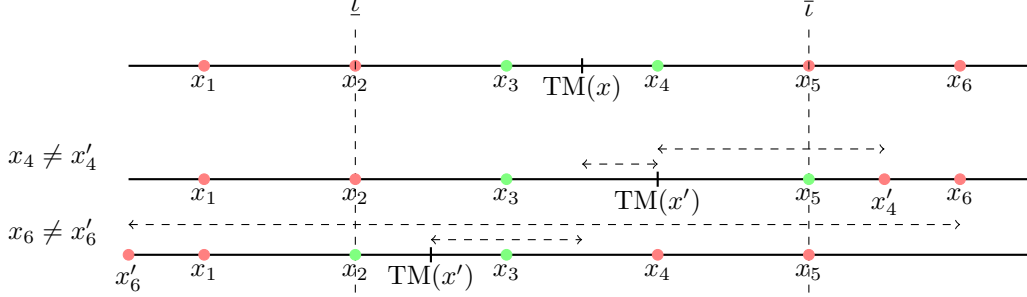


Figure 2: Example with $n = 6$, $f = 2$ and the differing point being either x_6 or x_4 . Horizontal dashed lines represent either trimmed mean differences or changed-point differences. Green points indicate the contributors to the trimmed-mean, while red points denote those excluded from the computation.

We consider the two possible cases.

- Case (i) $|S_{x'} \setminus S_x| = 1$. Let x_ℓ and $x_{\bar{\ell}}$ be the boundary elements not included in S_x . That is, if we denote by σ a permutation on $\{1, \dots, n\}$ that sorts the values in non-decreasing order $x_{\sigma(1)} \leq \dots \leq x_{\sigma(n)}$, then ℓ and $\bar{\ell}$ are defined as follow: $\sigma(\ell) = f$ and $\sigma(\bar{\ell}) = n - f + 1$.
 - If $\iota \notin S_x$ and $\iota \notin S_{x'}$, then with $\{j\} = S_x \setminus S_{x'}$, we have two sub-cases:
 - Sub-case (a) $x_\ell > x_{\bar{\ell}}$ and $x'_\ell < x_\ell$. In this case, $S_{x'} \setminus S_x = \{x_\ell\}$. Therefore,

$$|TM(x) - TM(x')| = \frac{1}{n-2f} |x_j - x_\ell| \leq \frac{1}{n-2f} |x_{\bar{\ell}} - x_\ell| \leq \frac{1}{n-2f} |x_\ell - x'_\ell|.$$
 - Sub-case (b) $x_\ell < x_{\bar{\ell}}$ and $x'_\ell > x_{\bar{\ell}}$. In this case, $S_{x'} \setminus S_x = \{x_{\bar{\ell}}\}$. Therefore,

$$|TM(x) - TM(x')| = \frac{1}{n-2f} |x_{\bar{\ell}} - x_j| \leq \frac{1}{n-2f} |x_{\bar{\ell}} - x_\ell| \leq \frac{1}{n-2f} |x_\ell - x'_\ell|.$$
 - If $\iota \in S_x$ but $\iota \notin S_{x'}$, then we have the following two sub-cases:
 - Sub-case (c) $x_\ell \in S_{x'}$. In this case, $x'_\ell \leq x_\ell$. Therefore,

$$|TM(x) - TM(x')| = \frac{1}{n-2f} |x_\ell - x_\ell| \leq \frac{1}{n-2f} |x'_\ell - x_\ell|.$$
 - Sub-case (d) $x_{\bar{\ell}} \in S_{x'}$. In this case, $x'_\ell \geq x_{\bar{\ell}}$. Therefore,

$$|TM(x) - TM(x')| = \frac{1}{n-2f} |x_{\bar{\ell}} - x_\ell| \leq \frac{1}{n-2f} |x'_\ell - x_\ell|.$$

- By symmetry, we obtain the same result for the case when $\iota \notin S_x$ but $\iota \in S_{x'}$.
- Case (ii) $|S_{x'} \setminus S_x| = 0$. The lemma is trivially true in this case.

From above, we obtain that $|\text{TM}(x) - \text{TM}(x')| \leq \frac{1}{n-2f} \|x - x'\|_1$.

We can generalize this bound to two sets of n real numbers that can differ in all the entries, using triangle inequality. Specifically,

$$\begin{aligned}
|\text{TM}(x) - \text{TM}(x')| &= |\text{TM}(x) - \text{TM}(x_1, \dots, x_{n-1}, x'_n) \\
&\quad + \text{TM}(x_1, \dots, x_{n-1}, x'_n) - \text{TM}(x_1, \dots, x_{n-2}, x'_{n-1}, x'_n) \\
&\quad + \dots + \text{TM}(x_1, x'_2, \dots, x'_n) - \text{TM}(x')| \\
&\leq |\text{TM}(x) - \text{TM}(x_1, \dots, x_{n-1}, x'_n)| \\
&\quad + \dots + |\text{TM}(x_1, x'_2, \dots, x'_n) - \text{TM}(x')| \\
&\leq \frac{1}{n-2f} \sum_{i=1}^n |x_i - x'_i| = \frac{1}{n-2f} \|x - x'\|_1 \leq \frac{\sqrt{n}}{n-2f} \|x - x'\|_2 \quad \square
\end{aligned}$$

We next present another property of the Trimmed Mean that, like Lipschitz continuity, enables us to establish the expansivity of the Component-Wise Trimmed Mean update rule.

Lemma D.7. *Let $x, x' \in \mathbb{R}^n$. There exist indices $\alpha, \beta \in \{1, \dots, n\}$ such that:*

$$[x - x']_\beta \leq \text{TM}(x) - \text{TM}(x') \leq [x - x']_\alpha.$$

Proof. Let $n \in \mathbb{N}$, $f \in \{0, \dots, \lfloor \frac{n}{2} \rfloor\}$, and $x, x' \in \mathbb{R}^n$

(i) Suppose $\text{TM}(x) = \text{TM}(x')$. We then consider two cases:

- If there exists an index $i \in \{1, \dots, n\}$ such that $x_i = x'_i$, the lemma is trivially true with $\alpha = \beta = i$.
- Otherwise, $x_i \neq x'_i, \forall i \in \{1, \dots, n\}$. In this case, there must be at least one index β where $[x - x']_\beta < 0$ and at least one index α where $[x - x']_\alpha > 0$. If all differences were strictly positive or strictly negative, this would contradict the assumption that $\text{TM}(x) = \text{TM}(x')$.

(ii) Suppose $\text{TM}(x) \neq \text{TM}(x')$. Here, we invoke the *translation equivariance* property of the trimmed mean. Specifically, let $\mathbf{1} = (1, \dots, 1) \in \mathbb{R}^n$, and define the translation operator $\mathcal{T}_{a\mathbf{1}}(x) := x - a\mathbf{1}$ for any $a \in \mathbb{R}$. Then, for any $a \in \mathbb{R}$, the trimmed mean satisfies:

$$\text{TM}(\mathcal{T}_{a\mathbf{1}}(x)) = \text{TM}(x - a\mathbf{1}) = \text{TM}(x) - a$$

Substituting $a = \text{TM}(x) - \text{TM}(x')$, we have $\text{TM}(\mathcal{T}_{a\mathbf{1}}(x)) - \text{TM}(x') = 0$. We then invoke the first case: there exist indices $\alpha, \beta \in \{1, \dots, n\}$ such that:

$$\begin{aligned}
[\mathcal{T}_{a\mathbf{1}}(x) - x']_\beta &\leq \text{TM}(\mathcal{T}_{a\mathbf{1}}(x)) - \text{TM}(x') \leq [\mathcal{T}_{a\mathbf{1}}(x) - x']_\alpha, \\
\implies [x - x']_\beta - a &\leq \text{TM}(x) - \text{TM}(x') - a \leq [x - x']_\alpha - a, \\
\implies [x - x']_\beta &\leq \text{TM}(x) - \text{TM}(x') \leq [x - x']_\alpha.
\end{aligned}$$

This concludes the proof. \square

Lemma D.6 and Lemma D.7 allow us to establish the expansivity of the Component-Wise Trimmed Mean update in the case of nonconvex loss functions, as shown in the following result.

Lemma D.8 (Expansivity of robust stochastic gradient update with CWTM). *Define the robust gradient update with learning rate $\gamma > 0$, for arbitrary $z^{(i)} \in \mathcal{Z}$, any $i \in \{1, \dots, n\}$, and any coordinate $k \in \{1, \dots, d\}$, as follows:*

$$[G_\gamma^{\text{CWTM}}(\theta)]_k = [\theta]_k - \gamma \text{TM}\left(\left[\nabla \ell(\theta; z^{(1)})\right]_k, \dots, \left[\nabla \ell(\theta; z^{(n)})\right]_k\right)$$

$$= [\theta]_k - \frac{\gamma}{n-2f} \sum_{i=f+1}^{n-f} [\nabla \ell(\theta; z^{(\sigma_k(i))})]_k = [\theta]_k - \frac{\gamma}{n-2f} \sum_{i \in S^{(k)}} [\nabla \ell(\theta; z^{(i)})]_k.$$

where σ_k denotes the permutation that sorts the values $([\nabla \ell(\theta; z^{(1)})]_k, \dots, [\nabla \ell(\theta; z^{(n)})]_k)$ in ascending order, and $S^{(k)} = \{i \in [n]; f+1 \leq \sigma_k(i) \leq n-f\}$ denotes the index set of central values in coordinate k . We can write,

$$G_\gamma^{\text{CWTM}}(\theta) = \theta - \gamma \text{CWTM}(\nabla \ell(\theta; z^{(1)}), \dots, \nabla \ell(\theta; z^{(n)})).$$

Assume $\forall z \in \mathcal{Z}, \ell(\cdot, z) : \Theta \rightarrow \mathbb{R}$ is L -smooth and nonconvex. Then G_γ^{CWTM} is $(1 + \gamma L \min\{\frac{n}{n-2f}, \sqrt{n}, \sqrt{d}\})$ -expansive.

Proof. Let $\theta, \omega \in \mathbb{R}^d$ and $z^{(1)}, \dots, z^{(n)} \in \mathcal{Z}$. For convenience, we occasionally adopt the shorthand notation: $\ell(\cdot; z^{(i)}) = \ell_i(\cdot)$. Using Lemma D.6 and the L -smoothness assumption, we obtain that

$$\begin{aligned} & \|\text{CWTM}(\nabla \ell_1(\theta), \dots, \nabla \ell_n(\theta)) - \text{CWTM}(\nabla \ell_1(\omega), \dots, \nabla \ell_n(\omega))\|_2^2 \\ &= \sum_{k=1}^d (\text{TM}([\nabla \ell_1(\theta)]_k, \dots, [\nabla \ell_n(\theta)]_k) - \text{TM}([\nabla \ell_1(\omega)]_k, \dots, [\nabla \ell_n(\omega)]_k))^2 \\ &\leq \frac{n}{(n-2f)^2} \sum_{k=1}^d \|([\nabla \ell_1(\theta)]_k, \dots, [\nabla \ell_n(\theta)]_k) - ([\nabla \ell_1(\omega)]_k, \dots, [\nabla \ell_n(\omega)]_k)\|_2^2 \\ &= \frac{n}{(n-2f)^2} \sum_{k=1}^d \sum_{i=1}^n ([\nabla \ell_i(\theta)]_k - [\nabla \ell_i(\omega)]_k)^2 = \frac{n}{(n-2f)^2} \sum_{i=1}^n \|\nabla \ell_i(\theta) - \nabla \ell_i(\omega)\|_2^2 \\ &\leq \frac{n}{(n-2f)^2} \sum_{i=1}^n L^2 \|\theta - \omega\|_2^2 = \frac{n^2 L^2}{(n-2f)^2} \|\theta - \omega\|_2^2. \end{aligned} \quad (11)$$

Alternatively, we can invoke Lemma D.7, which states that $\forall k \in \{1, \dots, d\}$, there exists $i_k \in \{1, \dots, n\}$ such that:

$$\begin{aligned} & \|\text{CWTM}(\nabla \ell_1(\theta), \dots, \nabla \ell_n(\theta)) - \text{CWTM}(\nabla \ell_1(\omega), \dots, \nabla \ell_n(\omega))\|_2^2 \\ &= \sum_{k=1}^d (\text{TM}([\nabla \ell_1(\theta)]_k, \dots, [\nabla \ell_n(\theta)]_k) - \text{TM}([\nabla \ell_1(\omega)]_k, \dots, [\nabla \ell_n(\omega)]_k))^2 \\ &\leq \sum_{k=1}^d [\nabla \ell_{i_k}(\theta) - \nabla \ell_{i_k}(\omega)]_k^2. \end{aligned}$$

From above, we obtain the following:

1. Since $\sum_{k=1}^d [\nabla \ell_{i_k}(\theta) - \nabla \ell_{i_k}(\omega)]_k^2 \leq \sum_{k=1}^d \sum_{i=1}^n [\nabla \ell_i(\theta) - \nabla \ell_i(\omega)]_k^2 \leq nL^2 \|\theta - \omega\|_2^2$,
 $\|\text{CWTM}(\nabla \ell_1(\theta), \dots, \nabla \ell_n(\theta)) - \text{CWTM}(\nabla \ell_1(\omega), \dots, \nabla \ell_n(\omega))\|_2^2 \leq nL^2 \|\theta - \omega\|_2^2$.
2. Since $\sum_{k=1}^d [\nabla \ell_{i_k}(\theta) - \nabla \ell_{i_k}(\omega)]_k^2 \leq \sum_{k=1}^d \|\nabla \ell_{i_k}(\theta) - \nabla \ell_{i_k}(\omega)\|_2^2$,
 $\|\text{CWTM}(\nabla \ell_1(\theta), \dots, \nabla \ell_n(\theta)) - \text{CWTM}(\nabla \ell_1(\omega), \dots, \nabla \ell_n(\omega))\|_2^2 \leq dL^2 \|\theta - \omega\|_2^2$.

Therefore,

$$\|\text{CWTM}(\nabla \ell_1(\theta), \dots, \nabla \ell_n(\theta)) - \text{CWTM}(\nabla \ell_1(\omega), \dots, \nabla \ell_n(\omega))\|_2^2 \leq \min\{n, d\} L^2 \|\theta - \omega\|_2^2. \quad (12)$$

The bound in (12) has a tighter dependence on the number of workers than the bound in (11), especially in the regime when f approaches $\frac{n}{2}$. We conclude the proof by combining (11) and (12), and by applying the triangle inequality, yielding:

$$\|G_\gamma^{\text{CWTM}}(\theta) - G_\gamma^{\text{CWTM}}(\omega)\|_2 \leq \left(1 + \gamma L \min\left\{\frac{n}{n-2f}, \sqrt{n}, \sqrt{d}\right\}\right) \|\theta - \omega\|_2. \quad \square$$

This behavior contrasts with the convex and strongly convex settings, where, as illustrated below in Remark D.3, a non-expansive update (*i.e.*, with expansivity coefficient ≤ 1) cannot, in general, be guaranteed. Notably, the non-expansivity property is retained in one-dimensional settings or when the loss function decomposes across dimensions, akin to the behavior observed with averaging-based updates. This observation is formalized in the following lemma.

Lemma D.9. *Let $d = 1$. Assume that for all $z \in \mathcal{Z}$, the loss function $\ell(\cdot; z)$ is convex and L -smooth. If $\gamma \leq \frac{2}{L}$ then G_γ^{CWTM} is 1-expansive.*

Proof. We use the same notation as in Lemma D.8, with $\ell'(\cdot)$ denoting the derivative of $\ell(\cdot)$. We obtain

$$\begin{aligned} |G_\gamma^{\text{CWTM}}(\theta) - G_\gamma^{\text{CWTM}}(\omega)|^2 &= |\theta - \omega|^2 + \gamma^2 |\text{TM}(\ell'_1(\theta), \dots, \ell'_n(\theta)) - \text{TM}(\ell'_1(\omega), \dots, \ell'_n(\omega))|^2 \\ &\quad - 2\gamma(\theta - \omega)(\text{TM}(\ell'_1(\theta), \dots, \ell'_n(\theta)) - \text{TM}(\ell'_1(\omega), \dots, \ell'_n(\omega))). \end{aligned} \quad (13)$$

In this case, we can leverage the co-coercivity property of the loss function. In fact, due to Lemma D.7, there exist indices $\alpha, \beta \in \{1, \dots, n\}$ such that:

$$\ell'_\beta(\theta) - \ell'_\beta(\omega) \leq \text{TM}(\ell'_1(\theta), \dots, \ell'_n(\theta)) - \text{TM}(\ell'_1(\omega), \dots, \ell'_n(\omega)) \leq \ell'_\alpha(\theta) - \ell'_\alpha(\omega). \quad (14)$$

Therefore, there exists $i \in \{\alpha, \beta\}$ such that

$$\begin{aligned} (\theta - \omega)(\text{TM}(\ell'_1(\theta), \dots, \ell'_n(\theta)) - \text{TM}(\ell'_1(\omega), \dots, \ell'_n(\omega))) &\geq (\theta - \omega)(\ell'_i(\theta) - \ell'_i(\omega)) \\ &\geq \frac{1}{L}(\ell'_i(\theta) - \ell'_i(\omega))^2 \geq 0. \end{aligned}$$

Thus, we have

$$\begin{aligned} &(\theta - \omega)(\text{TM}(\ell'_1(\theta), \dots, \ell'_n(\theta)) - \text{TM}(\ell'_1(\omega), \dots, \ell'_n(\omega))) \\ &= |\theta - \omega| |\text{TM}(\ell'_1(\theta), \dots, \ell'_n(\theta)) - \text{TM}(\ell'_1(\omega), \dots, \ell'_n(\omega))|. \end{aligned}$$

Using this in (13), we obtain that

$$\begin{aligned} |G_\gamma^{\text{CWTM}}(\theta) - G_\gamma^{\text{CWTM}}(\omega)|^2 &= |\theta - \omega|^2 + \gamma^2 |\text{TM}(\ell'_1(\theta), \dots, \ell'_n(\theta)) - \text{TM}(\ell'_1(\omega), \dots, \ell'_n(\omega))|^2 \\ &\quad - 2\gamma|\theta - \omega| |\text{TM}(\ell'_1(\theta), \dots, \ell'_n(\theta)) - \text{TM}(\ell'_1(\omega), \dots, \ell'_n(\omega))|. \end{aligned} \quad (15)$$

Note that, due to (14), there also exists $j \in \{\alpha, \beta\}$ such that

$$|\text{TM}(\ell'_1(\theta), \dots, \ell'_n(\theta)) - \text{TM}(\ell'_1(\omega), \dots, \ell'_n(\omega))| \leq |\ell'_j(\theta) - \ell'_j(\omega)| \leq L|\theta - \omega|. \quad (16)$$

Substituting from (16) in (15), we obtain that

$$\begin{aligned} |G_\gamma^{\text{CWTM}}(\theta) - G_\gamma^{\text{CWTM}}(\omega)|^2 &\leq |\theta - \omega|^2 \\ &\quad + \gamma(\gamma L - 2)|\theta - \omega| |\text{TM}(\ell'_1(\theta), \dots, \ell'_n(\theta)) - \text{TM}(\ell'_1(\omega), \dots, \ell'_n(\omega))|. \end{aligned}$$

The above proves the lemma since $(\gamma L - 2) \leq 0$ when $\gamma \leq \frac{2}{L}$. \square

Remark D.2. *Lemma D.9 generalizes to loss functions that separate across dimensions; that is, functions of the form $\ell = \sum_{k=1}^d \ell_k$, where each ℓ_k is a one dimensional function acting in the k -th coordinate, $k \in \{1, \dots, d\}$. A canonical example is mean estimation under coordinate-wise squared error loss, where the total loss decomposes as a sum of independent scalar losses for each dimension.*

Remark D.3. *To establish, in general, a non-expansivity result for the CWTM update rule, it would be necessary that the scalar product in:*

$$\begin{aligned} &\|G_\gamma^{\text{CWTM}}(\theta) - G_\gamma^{\text{CWTM}}(\omega)\|_2^2 \\ &= \|\theta - \omega\|_2^2 + \gamma^2 \|\text{CWTM}(\nabla \ell_1(\theta), \dots, \nabla \ell_n(\theta)) - \text{CWTM}(\nabla \ell_1(\omega), \dots, \nabla \ell_n(\omega))\|_2^2 \\ &\quad - 2\gamma \langle \text{CWTM}(\nabla \ell_1(\theta), \dots, \nabla \ell_n(\theta)) - \text{CWTM}(\nabla \ell_1(\omega), \dots, \nabla \ell_n(\omega)), \theta - \omega \rangle \end{aligned}$$

is strictly positive, that is the CWTM preserve the co-coercivity of the underlying loss function. As we will illustrate below, this is not the case. This is because the sets of selected indices differ across dimensions—that is $S^{(k)} \neq S^{(k)'} with $k \in \{1, \dots, d\}$ —which disrupts the alignment needed to preserve the co-coercivity property.$

Below, we present an example illustrating that CWTM does not preserve the co-coercivity property of an underlying convex and smooth loss function. Specifically, let $d = 2$, $n = 3$ and $f = 1$. We explicitly construct the example illustrated in Figure 3. We consider:

- the convex and smooth loss function:

$$\begin{aligned} \forall w \in \mathbb{R}^d, (x, y) \in B(0, \sqrt{L}) \times \mathbb{R} : \quad & \ell(w, (x, y)) = \frac{1}{2}(w^T x - y)^2 \\ L \in \mathbb{R}_+^* \quad & \nabla \ell(w, (x, y)) = (w^T x - y)x \\ & 0 \preceq \nabla^2 \ell(w, (x, y)) = xx^T \preceq LI_d \end{aligned}$$

- the sample pool composed by $z_1 = (v, 0)$, $z_2 = (x, 0)$ and $z_3 = (x, 1)$ for $x, v \in B(0, \sqrt{L})$.
- the parameters $\theta = \frac{x}{L}$ and $\omega = 0$.

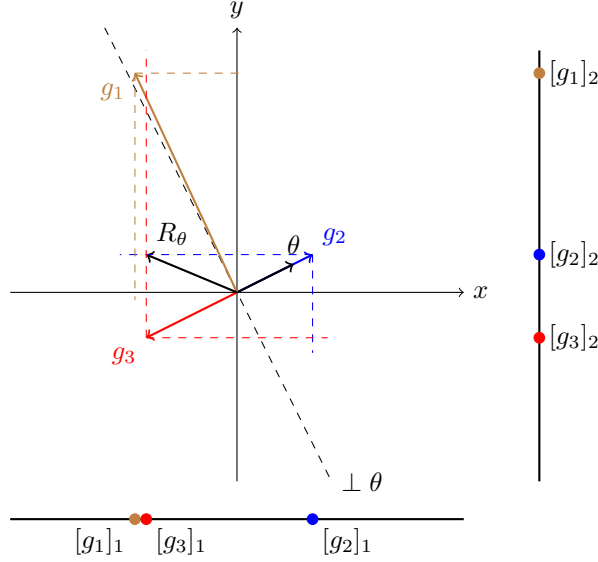


Figure 3: CWTM does not preserve the co-coercive inequality of smooth and convex functions.

We formally have: $g_1 = \theta^T v v$, $g_2 = \theta^T x x = \|x\|_2^2 \theta$ and $g_3 = \underbrace{(\|x\|_2^2 - L)}_{<0} \theta$. We next examine the conditions under which the CWTM aggregation rule does not preserve the co-coercivity inequality. Specifically, we denote: $R_\theta = \text{CWTM}(g_1, g_2, g_3)$ and $R_\omega = \text{CWTM}(0, 0, -x) = 0$:

$$\langle \theta - \omega, R_\theta - R_\omega \rangle = \langle \theta, R_\theta \rangle < 0$$

Instances of such v, x yielding this inequality are illustrated in the Figure 3 (v is colinear to g_1 and x colinear to g_2).

D.4 Proofs of Theorems 3.2 and 3.5

In the following section, we establish uniform algorithmic upper bounds under poisoning attacks, leveraging the structure of the SMEA robust aggregation rule or distinctive properties of the CWTM robust aggregation rule.

D.4.1 Proof of Theorem 3.2

We first prove upper bounds for robust distributed GD, and then a similar result for robust distributed SGD. We note that Theorem 3.2 merge the two similar results.

Theorem D.10 (Stability under poisoning attacks - SMEA – GD). *Consider the setting described in Section 2 under poisoning attacks. Let $\mathcal{A} = \text{GD}$, with SMEA, for $T \in \mathbb{N}^*$ iterations. Suppose $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ C -Lipschitz and L -smooth.*

(i) If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ convex, with constant learning rate $\gamma \leq \frac{2}{L}$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq 2\gamma C^2 T \left(\frac{f}{n-f} + \frac{1}{(n-f)m} \right).$$

(ii) If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ μ -strongly convex, with constant learning rate $\gamma \leq \frac{1}{L}$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq \frac{2C^2}{\mu} \left(\frac{f}{n-f} + \frac{1}{(n-f)m} \right).$$

Proof. Let denote $\{\theta_t\}_{t \in \{0, \dots, T-1\}}$ and $\{\theta'_t\}_{t \in \{0, \dots, T-1\}}$ the optimization trajectories resulting from two neighboring datasets $\mathcal{S}, \mathcal{S}'$, for T iterations of GD with aggregation rule SMEA and learning rate γ . For $t \in \{0, \dots, T-1\}$, let denote,

$$S_t^* \in \underset{S \subseteq \{1, \dots, n\}, |S|=n-f}{\operatorname{argmin}} \lambda_{\max} \left(\frac{1}{|S|} \sum_{i \in S} \left(g_t^{(i)} - \bar{g}_S \right) \left(g_t^{(i)} - \bar{g}_S \right)^\top \right)$$

with $\bar{g}_S = \frac{1}{|S|} \sum_{i \in S} g_t^{(i)}$, where $g_t^{(i)} = \nabla \hat{R}_i(\theta_t)$ and the prime notation $(\cdot)'$ denotes the corresponding quantities for \mathcal{S}' . At each step $t \in \{0, \dots, T-1\}$, the intersection and differences of the two subsets S_t^* and $S_t^{* \prime}$ verifies the following properties:

$$n - 2f \leq |S_t^* \cap S_t^{* \prime}| \leq n - f \quad \text{and} \quad |S_t^* \setminus S_t^{* \prime}| + |S_t^{* \prime} \setminus S_t^*| \leq 2f,$$

and let denote,

$$G_{\gamma_t | S_t^* \cap S_t^{* \prime}}^{\text{GD}}(\theta_t) = \theta_t - \frac{\gamma_t}{n-f} \sum_{i \in S_t^* \cap S_t^{* \prime}} \nabla \hat{R}_i(\theta_t),$$

and

$$G_{\gamma_t | S_t^* \cap S_t^{* \prime}}^{\text{GD}}(\theta'_t) = \theta'_t - \frac{\gamma_t}{n-f} \sum_{i \in S_t^* \cap S_t^{* \prime}} \nabla \hat{R}_i(\theta'_t).$$

Let $i \notin \mathcal{H}$, the misbehaving vectors are gradients computed on the true loss function. This enables us to exploit their regularities if $i \in S_t^* \cap S_t^{* \prime}$ or their boundedness if $i \in S_t^* \setminus S_t^{* \prime} \cup S_t^{* \prime} \setminus S_t^*$. This is to be contrasted with the previously studied Byzantine failure attacks—where an intermediate comparison was introduced—as we adopt a more direct bounding technique when subject to poisoning attacks. Using these relations, we derive the following bound on $\delta_{t+1} = \|\theta_{t+1} - \theta'_{t+1}\|_2$:

$$\begin{aligned} \delta_{t+1} &\leq \|G_{\gamma_t | S_t^* \cap S_t^{* \prime}}^{\text{GD}}(\theta_t) - G_{\gamma_t | S_t^* \cap S_t^{* \prime}}^{\text{GD}}(\theta'_t)\|_2 + \frac{2C\gamma_t}{(n-f)m} \\ &\quad + \frac{\gamma_t}{n-f} \left\| \sum_{i \in S_t^* \setminus S_t^{* \prime}} \nabla \hat{R}_i(\theta_t) - \sum_{i \in S_t^{* \prime} \setminus S_t^*} \nabla \hat{R}_i(\theta'_t) \right\|_2 \\ &\leq \|G_{\gamma_t | S_t^* \cap S_t^{* \prime}}^{\text{GD}}(\theta_t) - G_{\gamma_t | S_t^* \cap S_t^{* \prime}}^{\text{GD}}(\theta'_t)\|_2 + 2\gamma_t \frac{f + \frac{1}{m}}{n-f} C. \end{aligned}$$

Hence, by Lemma B.2, we obtain the following recursion:

$$\delta_{t+1} \leq \eta_{G_{\gamma_t}}^{\text{GD}} \delta_t + 2\gamma_t \frac{f + \frac{1}{m}}{n-f} C,$$

where we bound

$$\eta_{G_{\gamma_t | S_t^* \cap S_t^{* \prime}}^{\text{GD}}} = \frac{|S_t^* \cap S_t^{* \prime}|}{n-f} \eta_{G_{\gamma_t}}^{\text{GD}} \leq \eta_{G_{\gamma_t}}^{\text{GD}},$$

since $|S_t^* \cap S_t^{* \prime}| \leq n-f$. The precise value of the expansivity coefficient depends on the regularity assumption on the loss function. We then proceed the proof, as in Appendix D.2.1, depending on the regularities of the underlying loss function:

(i) In the Lipschitz, smooth, and convex case, we have for $t \in \{0, \dots, T-1\}$:

$$\delta_{t+1} \leq \delta_t + 2\gamma C \frac{f + \frac{1}{m}}{n-f}$$

We finish the proof by considering a constant learning rate $\gamma \leq \frac{1}{L}$ and by summing over $t \in \{0, \dots, T-1\}$ with $\theta_0 = \theta'_0$. We then invoke the Lipschitz continuity assumption.

(ii) In the Lipschitz, smooth and strongly convex case, we have for $t \in \{0, \dots, T-1\}$:

$$\delta_{t+1} \leq (1 - \gamma\mu)\delta_t + 2\gamma C \frac{f + \frac{1}{m}}{n-f}$$

Considering $\gamma \leq \frac{1}{L}$ and summing over $t \in \{0, \dots, T-1\}$ with $\theta_0 = \theta'_0$:

$$\delta_T \leq 2\gamma C \frac{f + \frac{1}{m}}{n-f} \sum_{t=0}^{T-1} (1 - \gamma\mu)^t = 2\gamma C \frac{f + \frac{1}{m}}{n-f} \frac{1}{\gamma\mu} [1 - (1 - \gamma\mu)^T] \leq \frac{2C}{\mu} \frac{f + \frac{1}{m}}{n-f}$$

We conclude the derivation of the bound by invoking the Lipschitz continuity assumption. \square

We now consider the SGD case, where for the sake of completeness we add a result for non-convex functions that was not presented in the main text.

Theorem D.11 (Stability under poisoning attacks - SMEA – SGD). *Consider the setting described in Section 2 under poisoning attacks. Let $\mathcal{A} = \text{SGD}$, with SMEA, for $T \in \mathbb{N}^*$ iterations. Suppose $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ C -Lipschitz and L -smooth.*

(i) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ convex, with constant learning rate $\gamma \leq \frac{2}{L}$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq 2\gamma C^2 T \left(\frac{f}{n-f} + \frac{1}{(n-f)m} \right).$$

(ii) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ μ -strongly convex, with constant learning rate $\gamma \leq \frac{1}{L}$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq \frac{2C^2}{\mu} \left(\frac{f}{n-f} + \frac{1}{(n-f)m} \right).$$

(iii) *If we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ nonconvex, bounded by $\ell_\infty = \sup_{\theta \in \Theta, z \in \mathcal{Z}} \ell(\theta; z)$, with a monotonically non-increasing learning rate $\gamma_t \leq \frac{c}{Lt}$, $t \in \{0, \dots, T-1\}$, $c > 0$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}}[|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq 2 \left(\frac{2C^2}{L} \frac{f + \frac{1}{m}}{n-f} \right)^{\frac{1}{c+1}} \left(\frac{\ell_\infty T}{m} \right)^{\frac{c}{c+1}}.$$

Proof. Let denote $\{\theta_t\}_{t \in \{0, \dots, T-1\}}$ and $\{\theta'_t\}_{t \in \{0, \dots, T-1\}}$ the optimization trajectories resulting from two neighboring datasets $\mathcal{S}, \mathcal{S}'$, for T iterations of SGD with aggregation rule SMEA and learning rate schedule $\{\gamma_t\}_{t \in \{0, \dots, T-1\}}$. The proof proceeds by analyzing our robust aggregation as a classical stochastic gradient descent, albeit with a random subset selection at each iteration. Notably, the poisoned vectors are treated as true gradients, although computed on arbitrary data. For $t \in \{0, \dots, T-1\}$, let denote,

$$z_t^{(i)} = \begin{cases} \text{an arbitrary poisoned data point,} & \text{if } i \notin \mathcal{H} \\ z^{(i, J_t^{(i)})}, J_t^{(i)} \text{ sampled uniformly from } \{1, \dots, m\} & \text{if } i \in \mathcal{H} \end{cases}$$

and

$$S_t^* \in \underset{S \subseteq \{1, \dots, n\}, |S|=n-f}{\operatorname{argmin}} \lambda_{\max} \left(\frac{1}{|S|} \sum_{i \in S} \left(g_t^{(i)} - \bar{g}_S \right) \left(g_t^{(i)} - \bar{g}_S \right)^\top \right)$$

with $\bar{g}_S = \frac{1}{|S|} \sum_{i \in S} g_t^{(i)}$ where $g_t^{(i)} = \nabla \ell(\theta_t; z_t^{(i)})$ and the prime notation $(\cdot)'$ denotes the corresponding quantities for S' . At each step $t \in \{t, \dots, T-1\}$, the intersection and differences of the two subsets S_t^* and $S_t'^*$ verifies the following properties:

$$n - 2f \leq |S_t^* \cap S_t'^*| \leq n - f \quad \text{and} \quad |S_t^* \setminus S_t'^*| + |S_t'^* \setminus S_t^*| \leq 2f,$$

and let denote,

$$G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta_t) = \theta_t - \frac{\gamma_t}{n-f} \sum_{i \in S_t^* \cap S_t'^*} g_t^{(i)},$$

and

$$G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta'_t) = \theta'_t - \frac{\gamma_t}{n-f} \sum_{i \in S_t^* \cap S_t'^*} g_t^{(i)'}$$

Let $i \notin \mathcal{H}$, the misbehaving vectors are gradients computed on the true loss function. This enables us to exploit their regularities if $i \in S_t^* \cap S_t'^*$ or their boundedness if $i \in S_t^* \setminus S_t'^* \cup S_t'^* \setminus S_t^*$. This is to be contrasted with the previously studied Byzantine failure attacks—where an intermediate comparison was introduced—as we adopt a more direct bounding technique when subject to poisoning attacks. Using these relations, we derive the following bound on $\delta_{t+1} = \|\theta_{t+1} - \theta'_{t+1}\|_2$, with probability $1 - \frac{1}{m}$:

$$\begin{aligned} \delta_{t+1} &\leq \|G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta_t) - G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta'_t)\|_2 \\ &\quad + \frac{\gamma_t}{n-f} \left\| \sum_{i \in S_t^* \setminus S_t'^*} \nabla \ell(\theta_t, z_t^{(i)}) - \sum_{i \in S_t'^* \setminus S_t^*} \nabla \ell(\theta'_t, z_t'^{(i)}) \right\|_2 \\ &\leq \|G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta_t) - G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta'_t)\|_2 + 2\gamma_t \frac{f}{n-f} C, \end{aligned}$$

and with probability $\frac{1}{m}$:

$$\begin{aligned} \delta_{t+1} &\leq \|G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta_t) - G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}'}(\theta'_t)\|_2 \\ &\quad + \frac{\gamma_t}{n-f} \left\| \sum_{i \in S_t^* \setminus S_t'^*} \nabla \ell(\theta_t, z_t^{(i)}) - \sum_{i \in S_t'^* \setminus S_t^*} \nabla \ell(\theta'_t, z_t'^{(i)}) \right\|_2 \\ &\leq \|G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta_t) - G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta'_t)\|_2 + \frac{\gamma_t}{n-f} \|\nabla \ell(\theta_t, z^{(a,b)}) - \nabla \ell(\theta'_t, z_t'^{(a,b)})\|_2 \\ &\quad + \frac{\gamma_t}{n-f} \left\| \sum_{i \in S_t^* \setminus S_t'^*} \nabla \ell(\theta_t, z_t^{(i)}) - \sum_{i \in S_t'^* \setminus S_t^*} \nabla \ell(\theta'_t, z_t'^{(i)}) \right\|_2 \\ &\leq \|G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta_t) - G_{\gamma_t | S_t^* \cap S_t'^*}^{\text{SGD}}(\theta'_t)\|_2 + 2\gamma_t \frac{f+1}{n-f} C, \end{aligned}$$

where the second inequality corresponds to the worst-case scenario, namely when the index of the worker with the differing data samples lies in the intersection of the two selected subsets of workers, i.e., $a \in S_t^* \cap S_t'^*$. Consequently, from Lemma B.2, we obtain the following recursion:

$$\mathbb{E}_{\mathcal{A}} \delta_{t+1} \leq \eta_{G_{\gamma_t}}^{\text{SGD}} \mathbb{E}_{\mathcal{A}} \delta_t + 2\gamma_t \frac{f + \frac{1}{m}}{n-f} C$$

where we bound

$$\eta_{G_{\gamma_t | S_t^* \cap S_t'^*}}^{\text{SGD}} = \frac{|S_t^* \cap S_t'^*|}{n-f} \eta_{G_{\gamma_t}}^{\text{SGD}} \leq \eta_{G_{\gamma_t}}^{\text{SGD}},$$

since $|S_t^* \cap S_t'^*| \leq n-f$. The precise value of the expansivity coefficient depends on the regularity assumption on the loss function. We then proceed the proof, as in Appendix D.2.1, depending on the regularities of the underlying loss function:

(i) In the Lipschitz, smooth and convex case, with $\gamma \leq \frac{2}{L}$, we have for $t \in \{0, \dots, T-1\}$:

$$\mathbb{E}_{\mathcal{A}} \delta_{t+1} \leq \mathbb{E}_{\mathcal{A}} \delta_t + 2\gamma \frac{f + \frac{1}{m}}{n-f} C$$

Unrolling the recursion and using the Lipschitz continuity of the loss function leads to:

$$\mathbb{E}_{\mathcal{A}} [\|\ell(\theta_T; z) - \ell(\theta'_T; z)\|] \leq 2\gamma C^2 T \frac{f + \frac{1}{m}}{n-f}$$

- (ii) In the Lipschitz, smooth and strongly convex case, with $\gamma \leq \frac{1}{L}$, we have for $t \in \{0, \dots, T-1\}$:

$$\mathbb{E}_{\mathcal{A}} \delta_{t+1} \leq (1 - \mu\gamma) \mathbb{E}_{\mathcal{A}} \delta_t + 2\gamma \frac{f + \frac{1}{m}}{n - f} C$$

Unrolling the recursion and using the Lipschitz continuity of the loss function yields:

$$\mathbb{E}_{\mathcal{A}} [|\ell(\theta_T; z) - \ell(\theta'_T; z)|] \leq 2\gamma C^2 \frac{f + \frac{1}{m}}{n - f} \frac{1}{\gamma\mu} \left[1 - (1 - \gamma\mu)^T\right] \leq \frac{2C^2}{\mu} \frac{f + \frac{1}{m}}{n - f}$$

- (iii) If the loss function is Lipschitz, smooth, and nonconvex, $c > 0$, $\gamma_t \leq \frac{c}{Lt}$, with $t_0 \in \{0, \dots, T-1\}$, for all $t \in \{t_0, \dots, T-1\}$, we obtain:

$$\mathbb{E}_{\mathcal{A}} [\delta_{t+1} | \delta_{t_0} = 0] \leq \left(1 + \frac{c}{t}\right) \mathbb{E}_{\mathcal{A}} [\delta_t | \delta_{t_0} = 0] + \frac{2c}{Lt} \frac{f + \frac{1}{m}}{n - f} C$$

Building on Lemma D.4, and following the same derivation as in Appendix D.2.2, we unroll the recursion, minimize over t_0 , and obtain:

$$\mathbb{E}_{\mathcal{A}} [|\ell(\theta_T; z) - \ell(\theta'_T; z)|] \leq 2 \left(\frac{2C^2}{L} \frac{f + \frac{1}{m}}{n - f} \right)^{\frac{1}{c+1}} \left(\frac{\ell_{\infty} T}{m} \right)^{\frac{c}{c+1}}. \quad \square$$

Remark D.4. Our bounds does not depend on the robustness coefficient κ of the SMEA robust aggregation, but solely on its structural characteristic—namely, the sub-sampling of worker vectors. Additionally, our proof does not exploit the randomness in the selection of workers at each step, as this would require assumptions regarding the local distributions of the misbehaving workers to evaluate the probability of selecting the differing data samples' gradients. Interestingly, in practice, one could observe improved stability compared to the setting without misbehaving workers, possibly due to an extended “burn-in period”. That is, if the differing samples are not selected by the aggregation rule when first drawn. This phenomenon could be further explored through numerical experiments.

Remark D.5. We can actually achieve the same theoretical bound as the non-misbehaving workers framework without the robust aggregation rule. Specifically, by reverting to the classical non-robust mean aggregation, we are able to handle terms as described in [23], since data heterogeneity does not manifest in the worst-case uniform algorithmic stability analysis. However, it is important to note that while this approach simplifies the analysis, the population risk may still remain unbounded due to the inherently unbounded optimization error that can arise from mean aggregation when facing an arbitrarily poisoned dataset.

D.4.2 Proof of Theorem 3.5

In the following, we derive sharper uniform stability upper bound with CWTM aggregation rule for smooth and non-convex loss functions. This improvement—particularly in its dependence on the number of local samples—is enabled by aggregation rules that preserve the smoothness of the loss function

Theorem 3.5 (Stability under poisoning attacks - CWTM – SGD). *Consider the setting described in Section 2 under poisoning attacks. Let $\mathcal{A} = \text{SGD}$, with CWTM. Suppose \mathcal{A} is run for $T \in \mathbb{N}^*$ iterations and that there exist $\nu > 0$ constant such that $n \geq (2 + \nu)f$ — that is, f/n is strongly bounded away from $1/2$. Let we assume $\forall z \in \mathcal{Z}$, $\ell(\cdot; z)$ bounded by ℓ_{∞} , nonconvex, C -Lipschitz and L -smooth. Then, with a monotonically non-increasing learning rate, $\gamma_t \leq \frac{\nu}{2+\nu} \frac{c}{Lt}$, $t \in \{0, \dots, T-1\}$, $c > 0$, we have for any neighboring datasets $\mathcal{S}, \mathcal{S}'$:*

$$\sup_{z \in \mathcal{Z}} \mathbb{E}_{\mathcal{A}} [|\ell(\mathcal{A}(\mathcal{S}); z) - \ell(\mathcal{A}(\mathcal{S}'); z)|] \leq 2 \left(\frac{2C^2 \nu^2}{(2 + \nu)^2 L} \right)^{\frac{1}{c+1}} \frac{(T \ell_{\infty})^{\frac{c}{c+1}}}{m \sqrt{n}^{\frac{1}{c+1}}}$$

Proof. Let denote $\{\theta_t\}_{t \in \{0, \dots, T-1\}}$ and $\{\theta'_t\}_{t \in \{0, \dots, T-1\}}$ the optimization trajectories resulting from two neighboring datasets $\mathcal{S}, \mathcal{S}'$, for T iterations of SGD with aggregation rule CWTM and learning rate schedule $\{\gamma_t\}_{t \in \{0, \dots, T-1\}}$. Importantly, poisoned vectors are treated as legitimate gradients, despite being computed on arbitrary data. In the following, we draw a growth recursion for the

sensitivity of the parameters to the perturbation introduced. Let $a, b \in \mathcal{H} \times \{1, \dots, m\}$ the indices of the differing samples. For $t \in \{0, \dots, T-1\}$, we denote:

$$z_t^{(i)} = \begin{cases} \text{an arbitrary poisoned data point,} & \text{if } i \notin \mathcal{H} \\ z^{(i, J_t^{(i)})}, J_t^{(i)} \text{ sampled uniformly from } \{1, \dots, m\} & \text{if } i \in \mathcal{H} \end{cases}$$

$\delta_t = \|\theta_t - \theta'_t\|_2$, $g_t^{(i)} = \nabla \ell(\theta_t, z_t^{(i)})$, $G_{\gamma_t}^{\text{CWTM}}(\theta_t) = \theta_t - \gamma_t \text{CWTM}(g_t^{(1)}, \dots, g_t^{(n)})$, and the prime notation $(\cdot)'$ denotes the corresponding quantities for \mathcal{S}' .

The proof leverages the expansivity property of the robust aggregation update without directly comparing it to the honest stochastic gradient update. Specifically, we either directly apply the expansivity result Lemma D.8 with probability $(1 - \frac{1}{m})$:

$$\|G_{\gamma}^{\text{CWTM}}(\theta_t) - G_{\gamma}^{\text{CWTM}}(\theta'_t)\|_2 \leq (1 + \gamma_t L \min\{\frac{n}{n-2f}, \sqrt{n}, \sqrt{d}\}) \|\theta_t - \theta'_t\|_2$$

or perform the following calculation with probability $\frac{1}{m}$:

$$\begin{aligned} & \|\text{CWTM}(g_t^{(1)}, \dots, g_t^{(n)}) - \text{CWTM}(g_t'^{(1)}, \dots, g_t'^{(n)})\|_2^2 \\ &= \sum_{k=1}^d \left(TM([g_t^{(1)}]_k, \dots, [g_t^{(n)}]_k) - TM([g_t'^{(1)}]_k, \dots, [g_t'^{(n)}]_k) \right)^2 \\ &\leq \frac{n}{(n-2f)^2} \sum_{k=1}^d \left\| ([g_t^{(1)}]_k, \dots, [g_t^{(n)}]_k) - ([g_t'^{(1)}]_k, \dots, [g_t'^{(n)}]_k) \right\|_2^2 \\ &= \frac{n}{(n-2f)^2} \sum_{k=1}^d \sum_{i=1}^n \left([g_t^{(i)}]_k - [g_t'^{(i)}]_k \right)^2 \\ &= \frac{n}{(n-2f)^2} \sum_{i=1}^n \|\nabla \ell(\theta_t; z_t^{(i)}) - \nabla \ell(\theta'_t; z_t'^{(i)})\|_2^2 \\ &= \frac{n}{(n-2f)^2} \left(\sum_{i=1, i \neq a}^n \|\nabla \ell(\theta_t; z_t^{(i)}) - \nabla \ell(\theta'_t; z_t^{(i)})\|_2^2 \right. \\ &\quad \left. + \|\nabla \ell(\theta_t; z^{(a,b)}) - \nabla \ell(\theta'_t; z'^{(a,b)})\|_2^2 \right) \\ &\leq \frac{n}{(n-2f)^2} (4C^2 + (n-1)L^2 \|\theta_t - \theta'_t\|_2^2) \end{aligned}$$

where we used Lemma D.6 in the first inequality. Using the sub-additivity of the square root we prove that we have with probability $\frac{1}{m}$:

$$\|G_{\gamma}^{\text{CWTM}}(\theta_t) - G_{\gamma}^{\text{CWTM}}(\theta'_t)\|_2 \leq (1 + \frac{\gamma_t L \sqrt{n(n-1)}}{n-2f}) \|\theta_t - \theta'_t\|_2 + \frac{2\gamma_t C \sqrt{n}}{n-2f}.$$

Consequently, with $t_0 \in \{0, \dots, T-1\}$, we obtain the following recursion:

$$\begin{aligned} \mathbb{E}_{\mathcal{A}}[\delta_{t+1} | \delta_{t_0} = 0] &\leq (1 + \frac{\gamma_t L n}{n-2f}) \mathbb{E}_{\mathcal{A}}[\delta_t | \delta_{t_0} = 0] + \frac{2\gamma_t C \sqrt{n}}{(n-2f)m} \\ &\leq (1 + \frac{\gamma_t L (2+\nu)}{\nu}) \mathbb{E}_{\mathcal{A}}[\delta_t | \delta_{t_0} = 0] + \sigma_{f,m,n} \frac{c}{Lt} \\ &= (1 + \frac{c}{t}) \mathbb{E}_{\mathcal{A}}[\delta_t | \delta_{t_0} = 0] + \sigma_{f,m,n} \frac{c}{Lt}. \end{aligned}$$

where we used that as $n \geq (2+\nu)f$, then $\frac{n}{n-2f} \leq \frac{2+\nu}{\nu}$. and defined $\sigma_{f,m,n} = \frac{2C\nu\sqrt{n}}{(2+\nu)(n-2f)m} \leq \frac{2C\nu^2}{(2+\nu)^2\sqrt{nm}}$. We then derive our bound similarly to the case outlined in Appendix D.2.2. For

$0 \leq t_0 \leq \min(m, T)$, $t \in \{t_0, \dots, T-1\}$, $\gamma_t \leq \frac{\nu}{2+\nu} \frac{c}{Lt}$, $c > 0$, we use $1+x \leq e^x$ and unroll the recursion:

$$\mathbb{E}_{\mathcal{A}}[\delta_T | \delta_{t_0} = 0] \leq \frac{c\sigma_{f,m,n}}{L} \sum_{t=t_0}^{T-1} \frac{1}{t} \prod_{s=t+1}^{T-1} e^{\frac{c}{s}} \leq \frac{\sigma_{f,m,n}}{L} \left(\frac{T}{t_0}\right)^c$$

because:

$$\begin{aligned} \sum_{t=t_0}^{T-1} \frac{1}{t} \prod_{s=t+1}^{T-1} e^{\frac{c}{s}} &= \sum_{t=t_0}^{T-1} \frac{1}{t} e^{c \sum_{s=t+1}^{T-1} \frac{1}{s}} \leq \sum_{t=t_0}^{T-1} \frac{1}{t} e^{c \log(\frac{T}{t})} \\ &= T^c \sum_{t=t_0}^{T-1} t^{-c-1} \stackrel{\text{see Appendix D.2.2}}{\leq} \frac{1}{c} \left(\frac{T}{t_0}\right)^c \end{aligned}$$

Finally, considering Lemma D.4, we obtain:

$$\begin{aligned} \mathbb{E}_{\mathcal{A}}[|\ell(\theta_T; z) - \ell(\theta'_T; z)|] &\leq \min_{0 \leq t_0 \leq \min(m, T)} \frac{t_0 \ell_\infty}{m} + \frac{\sigma_{f,m,n} C}{L} \left(\frac{T}{t_0}\right)^c \\ &\leq \min_{0 \leq t_0 \leq \min(m, T)} \frac{t_0 \ell_\infty}{m} + \frac{2C^2 \nu^2}{(2+\nu)^2 L \sqrt{nm}} \left(\frac{T}{t_0}\right)^c. \end{aligned}$$

We approximately minimize the expression with respect to t_0 by balancing the two terms, *i.e.*, setting them equal:

$$0 \leq \tilde{t}_0 = \left(\frac{m\sigma_{f,m,n}C}{L\ell_\infty}\right)^{\frac{1}{c+1}} T^{\frac{c}{c+1}} \leq T \quad \text{for sufficiently large } T \geq \frac{m\sigma_{f,m,n}C}{L\ell_\infty}.$$

Finally, we obtain the result by directly substituting \tilde{t}_0 into the bound:

$$\begin{aligned} \mathbb{E}_{\mathcal{A}}[|\ell(\theta_T; z) - \ell(\theta'_T; z)|] &\leq 2 \left(\frac{\sigma_{f,m,n}C}{L}\right)^{\frac{1}{c+1}} \left(\frac{\ell_\infty T}{m}\right)^{\frac{c}{c+1}} \\ &= 2 \left(\frac{2C^2 \nu^2}{(2+\nu)^2 L \sqrt{nm}}\right)^{\frac{1}{c+1}} \left(\frac{\ell_\infty T}{m}\right)^{\frac{c}{c+1}} = 2 \left(\frac{2C^2 \nu^2}{(2+\nu)^2 L}\right)^{\frac{1}{c+1}} (\ell_\infty T)^{\frac{c}{c+1}} \frac{1}{m \sqrt{n}^{\frac{1}{c+1}}} \quad \square \end{aligned}$$

D.4.3 Comparaision of upper bounds in nonconvex settings

In this section, we present additional comparisons of upper bounds in nonconvex settings that were not included in the main text. Specifically, we compare the bounds obtained using the SMEA aggregation rule under poisoning attacks with those obtained using CWTM, as well as with bounds under Byzantine failure attacks.

- (i) We compare $\varepsilon^{\text{Byzantine}}$ from eq. (6) to $\varepsilon_{\text{SMEA}}^{\text{poisoning}} = 2 \left(\frac{2C^2}{L} \frac{f+\frac{1}{m}}{n-f}\right)^{\frac{1}{c+1}} \left(\frac{\ell_\infty T}{m}\right)^{\frac{c}{c+1}}$ from Theorem D.11: $\varepsilon^{\text{Byzantine}} / \varepsilon_{\text{SMEA}}^{\text{poisoning}} = \left(\frac{\frac{1}{(n-f)m} + \sqrt{\kappa}}{\frac{1}{(n-f)m} + \frac{f}{n-f}}\right)^{\frac{1}{c+1}}$. While the dependence on m remains unchanged, we demonstrate an improvement from $\sqrt{\kappa}$ to $\frac{f}{n-f}$.
- (ii) We compare $\varepsilon_{\text{SMEA}}^{\text{poisoning}}$ from theorem D.11 to $\varepsilon_{\text{CWTM}}^{\text{poisoning}}$ from eq. (7): $\varepsilon_{\text{SMEA}}^{\text{poisoning}} / \varepsilon_{\text{CWTM}}^{\text{poisoning}} = \left(\frac{(2+\nu)^2}{\nu^2} \left(\frac{\sqrt{n}}{n-f} + fm \frac{\sqrt{n}}{n-f}\right)\right)^{\frac{1}{c+1}}$. We observe a more favorable dependence on m in the CWTM bound.

D.5 Proof of Theorem 3.3 for GD

In this section, we derive uniform stability lower-bound for distributed GD with SMEA aggregation rule under poisoning attacks.

The following lemma plays a key role in understanding the SMEA aggregation rule, which is essential for the subsequent lower-bound constructions.

Lemma D.12. Let $n, f \in \mathbb{N}^*$, $f < \frac{n}{2}$ and g_A, g_B, g_C three collinear vectors in \mathbb{R}^d , $d \geq 1$. Assume we observe a pool of n vectors, consisting of α copies of g_A , β copies of g_B and γ copies of g_C , where $\alpha, \beta, \gamma \in \{0, \dots, n\}$, $\alpha + \beta + \gamma = n$. We abstract each subset $S \subseteq \{1, \dots, n\}$, composed of $|S| = n - f$ vectors from our pool of vectors as $S_{a,b,c}$ with $a \in \{0, \dots, \alpha\}$, $b \in \{0, \dots, \beta\}$, $c \in \{0, \dots, \gamma\}$, such that $a + b + c = n - f$ and S is composed by a copies of g_A , b copies of g_B and c copies of g_C . With this notation, we have:

$$\begin{aligned} S^* &\in \underset{\substack{S \subseteq \{1, \dots, n\}, \\ |S| = n-f}}{\operatorname{argmin}} \lambda_{\max} \left(\frac{1}{|S|} \sum_{i \in S} (g_i - \bar{g}_S)(g_i - \bar{g}_S)^\top \right) \\ &= \frac{1}{(n-f)^2} \underset{\substack{S_{a,b,c} \\ a \in \{0, \dots, \alpha\}, b \in \{0, \dots, \beta\}, \\ c \in \{0, \dots, \gamma\}, a+b+c=n-f}}{\operatorname{argmin}} \{ab\|g_A - g_B\|_2^2 + bc\|g_B - g_C\|_2^2 + ac\|g_A - g_C\|_2^2\} \end{aligned}$$

Let u be a unit vector indicating the line the vectors are colinear to. Since all computations are effectively confined to this one-dimensional subspace, we could equivalently emphasize that:

$$\lambda_{\max} \left(\frac{1}{|S|} \sum_{i \in S} (g_i - \bar{g}_S)(g_i - \bar{g}_S)^\top \right) = \frac{1}{|S|} \sum_{i \in S} (g_i^\top u - \bar{g}_S^\top u)^2 \lambda_{\max}(uu^\top) = \operatorname{Var}(\{g_i^\top u\}_{i \in S})$$

Proof. Let $S_{a,b,c}$ with $a \in \{0, \dots, \alpha\}$, $b \in \{0, \dots, \beta\}$, $c \in \{0, \dots, \gamma\}$, such that $a + b + c = n - f$ and S is composed by a copies of g_A , b copies of g_B and c copies of g_C . Then $\bar{g}_S = \frac{1}{n-f}(ag_A + bg_B + cg_C)$, $g_A - \bar{g}_S = \frac{b}{n-f}(g_A - g_B) + \frac{c}{n-f}(g_A - g_C)$, $g_B - \bar{g}_S = \frac{a}{n-f}(g_B - g_A) + \frac{c}{n-f}(g_B - g_C)$ and $g_C - \bar{g}_S = \frac{a}{n-f}(g_C - g_A) + \frac{b}{n-f}(g_C - g_B)$. We already note that the outer product of vectors g_A , g_B and g_C commute because they are collinear. We denote $(*) = (n-f)^2 \sum_{i \in S_{a,b,c}} (g_i - \bar{g}_S)(g_i - \bar{g}_S)^\top$, by expanding and re-arranging the terms:

$$\begin{aligned} (*) &= a \times [b^2(g_A - g_B)(g_A - g_B)^\top + c^2(g_A - g_C)(g_A - g_C)^\top + 2bc(g_A - g_B)(g_A - g_C)^\top] \\ &\quad + b \times [a^2(g_B - g_A)(g_B - g_A)^\top + c^2(g_B - g_C)(g_B - g_C)^\top + 2ac(g_B - g_A)(g_B - g_C)^\top] \\ &\quad + c \times [a^2(g_C - g_A)(g_C - g_A)^\top + b^2(g_C - g_B)(g_C - g_B)^\top + 2ab(g_C - g_A)(g_C - g_B)^\top] \\ &= (n-f)ab(g_A - g_B)(g_A - g_B)^\top - abc(g_A - g_B)(g_A - g_B)^\top \\ &\quad + (n-f)ac(g_A - g_C)(g_A - g_C)^\top - abc(g_A - g_C)(g_A - g_C)^\top \\ &\quad + (n-f)bc(g_B - g_C)(g_B - g_C)^\top - abc(g_B - g_C)(g_B - g_C)^\top \\ &\quad + 2abc[(g_A - g_B)(g_A - g_C)^\top + (g_B - g_A)(g_B - g_C)^\top + (g_C - g_A)(g_C - g_B)^\top] \end{aligned}$$

We then re-arrange the terms and simplify the result:

$$\begin{aligned} (*) &= (n-f)ab(g_A - g_B)(g_A - g_B)^\top - abc(g_A - g_B)(g_A - g_B)^\top \\ &\quad + (n-f)ac(g_A - g_C)(g_A - g_C)^\top - abc(g_A - g_C)(g_A - g_C)^\top \\ &\quad + (n-f)bc(g_B - g_C)(g_B - g_C)^\top - abc(g_B - g_C)(g_B - g_C)^\top \\ &\quad + abc[(g_A - g_B)(g_A - g_C)^\top - (g_A - g_B)(g_B - g_C)^\top] \\ &\quad + abc[(g_A - g_B)(g_A - g_C)^\top - (g_A - g_C)(g_C - g_B)^\top] \\ &\quad + abc[(g_B - g_A)(g_B - g_C)^\top - (g_C - g_A)(g_B - g_C)^\top] \\ &= (n-f)[ab(g_A - g_B)(g_A - g_B)^\top + ac(g_A - g_C)(g_A - g_C)^\top + bc(g_B - g_C)(g_B - g_C)^\top] \end{aligned}$$

To conclude the proof, we leverage the fact that g_A, g_B, g_C are collinear vectors. Let u be a unit vector representing the line the vectors are colinear to. Since each $g \in \{g_A, g_B, g_C\}$ lies along the line spanned by u , we can factor out u and use the fact that $\lambda_{\max}(uu^\top) = 1$:

$$\begin{aligned} \lambda_{\max} \left(\frac{1}{n-f} \sum_{i \in S_{a,b,c}} (g_i - \bar{g}_S)(g_i - \bar{g}_S)^\top \right) &= \frac{1}{(n-f)^2} \left(ab\|g_A - g_B\|_2^2 + bc\|g_B - g_C\|_2^2 \right. \\ &\quad \left. + ac\|g_A - g_C\|_2^2 \right) \quad \square \end{aligned}$$

D.5.1 Convex case

Theorem D.13 (Added instability in robust GD - convex case). *Consider the setting described in Section 2 under data poisoning attacks. Let $\mathcal{F}_{\mathcal{Z}}$ denote the set of functions that are C -Lipschitz, L -smooth and convex on \mathcal{Z} . Let $\mathcal{A} = \text{GD}$, with the SMEA. Suppose \mathcal{A} is run for $T \in \mathbb{N}^*$ iterations with $\gamma \leq \frac{2}{L}$. Then there exists a loss function $\ell \in \mathcal{F}_{\mathcal{Z}}$ and a pair of neighboring datasets such that the uniform stability is lower bounded by*

$$\Omega \left(\gamma C^2 T \left(\frac{f}{n-f} + \frac{1}{(n-f)m} \right) \right).$$

Proof. In what follows, we provide the proof for the case $d = 1$. The result extends naturally to higher dimensions by embedding the one-dimensional construction into \mathbb{R}^d . An explicit construction for $d = 2$ is presented in the lower-bound proof for SGD, with SMEA, and convex loss function. We arbitrarily refer to one trajectory as *unperturbed*, with parameter θ_t , and the other as *perturbed*, with θ'_t . The proof is structured in the following three parts:

- (i) First, we specify the loss function and the pair of neighboring datasets used. We then describe the structure of our problem, along with its regularity. Broadly speaking, the setup involves three distinct subgroups of workers. The first subgroup acts as a neutral or pivot subgroup that contains the differing samples, while the other two subgroups compete to be selected by the aggregation rule alongside the neutral subgroup.
- (ii) We then show that, within our setup, it is possible to favor one subgroup through initialization and the other through perturbation introduced by the differing sample. This property holds inductively throughout the optimization process, ensuring that the unperturbed and perturbed parameters converge to distinct minimizers. As a result, the two trajectories remain divergent at each iteration, introducing an additional expansive term relative to the behavior observed under standard averaging rule.
- (iii) Finally, we leverage this construction to derive a lower bound on uniform stability, which matches the upper bound established in Theorem D.10.

(i) Let $C, L \in \mathbb{R}_+^*$, we define the following loss function:

$$\text{For } \theta \in \mathbb{R}, z \in [-C, C], \quad \ell(\theta; z) = z\theta, \quad \ell'(\theta; z) = z, \quad \ell''(\theta; z) = 0. \quad (17)$$

By construction, the loss function is convex, C -Lipschitz and L -smooth. Let $m, n, f \in \mathbb{N}^*$ such that $f < \frac{n}{2}$. We define:

$$0 < \psi := \sqrt{\frac{n-2f-\frac{2}{m}}{n-2f+\frac{2}{m}}} < 1,$$

and the following pair of neighboring datasets:

$$\begin{aligned} S &= \{D_1, \dots, D_n\}, \quad S' = S \setminus \{D_1\} \cup \{D_1 \setminus \{z^{(1,1)}\} \cup \{z'^{(1,1)}\}\} \\ \text{with } \forall i \in \{1, \dots, n\}, \quad D_i &= \{z^{(i,1)}, \dots, z^{(i,m)}\} \\ \text{and } \forall j \in \{1, \dots, m\}, \quad i = 1: z^{(1,j)} &= 0, \quad z^{(1,1)} = -C, \text{ and } z'^{(1,1)} = 0. \\ i \in N &= \{2, \dots, n-2f\}, |N| = n-2f-1: z^{(i,j)} = 0. \\ i \in E &= \{n-2f+1, \dots, n-f\}, |E| = f: z^{(i,j)} = \frac{1+\psi}{2}C. \\ i \in F &= \{n-f+1, \dots, n\}, |F| = f: z^{(i,j)} = -C. \end{aligned} \quad (18)$$

Building on this setup, we first characterize the form of the parameters and gradients throughout the training process. We then specify the conditions under which our lower bound holds. Let $\gamma \leq \frac{2}{L}$ the learning rate and assume $\theta_0 = \theta'_0 = 0$. For $t \in \{0, \dots, T-1\}$, $g_t^{(N)} = g_t^{(N)'} = g_t^{(1)'} = 0$, $g_t^{(F)} = g_t^{(F)'} = -C$, $g_t^{(E)} = g_t^{(E)'} = \frac{1+\psi}{2}C$, and $g_t^{(1)} = -\frac{C}{m}$.

(ii) In what follows, we denote, for a subset $S \subset \{1, \dots, n\}$ of size $n - f$, and for all $t \in \{0, \dots, T - 1\}$,

$$\lambda_{\max}^S := \lambda_{\max} \left(\frac{1}{|S|} \sum_{i \in S} (g_t^{(i)} - \bar{g}_S)(g_t^{(i)} - \bar{g}_S)^\top \right) = \lambda_{\max} \left(\frac{1}{|S|} \sum_{i \in S} (g_0^{(i)} - \bar{g}_S)(g_0^{(i)} - \bar{g}_S)^\top \right)$$

with $\bar{g}_S = \frac{1}{|S|} \sum_{i \in S} g_0^{(i)}$. Recall that SMEA performs averaging over the subset of gradients selected according to the following rule:

$$S_t^* \in \underset{S \subseteq \{1, \dots, n\}, |S|=n-f}{\operatorname{argmin}} \lambda_{\max} \left(\frac{1}{|S|} \sum_{i \in S} (g_t^{(i)} - \bar{g}_S)(g_t^{(i)} - \bar{g}_S)^\top \right).$$

We highlight that the gradients are independent of the parameters and remain constant throughout the optimization process, as they are equal to the mean of their local samples.

With a slight abuse of notation, we denote, for instance, the set $\{1\} \cup N \cup E$ by $\{1, N, E\}$. To uncover the source of instability, we consider the following two key conditions:

(C1) By Lemma D.12, for all $t \in \{0, \dots, T - 1\}$, $S_t^{*'} = \{1, N, E\}$ is equivalent to

$$\lambda_{\max}^{\{1, N, F\}} = \frac{f(n - 2f)}{(n - f)^2} C^2 > \frac{f(n - 2f)}{(n - f)^2} \left(\frac{1 + \psi}{2} \right)^2 C^2 = \lambda_{\max}^{\{1, N, E\}},$$

which is ensured by construction. In fact, we again recall that the gradients are independent of the parameters and remain constant throughout the optimization process. Next, we demonstrate that any subset composed of a mixture of gradients from 1, N , F , and E necessarily yields a larger maximum eigenvalue. This is due to the fact that the vectors $g_0^{(F)'} and $g_0^{(E)'}$ are collinear but oriented in opposite directions, thereby contributing additively to the overall magnitude. To make this precise, we consider the following two cases:$

(1) Let denote S the set were we remove $q < f$ gradients of the subgroup E and add q gradients of the subgroup F from the set $\{1, N, E\}$, we end up comparing the previous quantity to the following, by Lemma D.12:

$$\begin{aligned} \lambda_{\max}^{S'} &= \frac{(f - q)(n - 2f)}{(n - f)^2} C^2 \|g_0^{(E)'} - g_0^{(N)'}\|_2^2 + \frac{q(n - 2f)}{(n - f)^2} \|g_0^{(F)'} - g_0^{(N)'}\|_2^2 \\ &\quad + \frac{q(f - q)}{(n - f)^2} \|g_0^{(F)'} - g_0^{(E)'}\|_2^2 \\ &= \frac{(f - q)(n - 2f)}{(n - f)^2} C^2 \left(\frac{1 + \psi}{2} \right)^2 + \frac{q(n - 2f)}{(n - f)^2} C^2 + q(f - q) C^2 \left(\frac{3 + \psi}{2} \right)^2 \\ &= \lambda_{\max}^{\{1, N, E\}} + \frac{q(n - 2f)}{(n - f)^2} C^2 \left(1 - \left(\frac{1 + \psi}{2} \right)^2 \right) + \frac{q(f - q)}{(n - f)^2} C^2 \left(\frac{3 + \psi}{2} \right)^2 \\ &> \lambda_{\max}^{\{1, N, E\}}. \end{aligned}$$

(2) Similarly, let denote S the set were we remove $q < \min(f, n - 2f)$ gradients of the subgroup $\{1, N\}$ and add q gradients of the subgroup F from the set $\{1, N, E\}$, we end up comparing the previous quantity to the following, by Lemma D.12:

$$\begin{aligned} \lambda_{\max}^{S'} &= \frac{f(n - 2f - q)}{(n - f)^2} C^2 \|g_0^{(E)'} - g_0^{(N)'}\|_2^2 + \frac{q(n - 2f - q)}{(n - f)^2} \|g_0^{(F)'} - g_0^{(N)'}\|_2^2 \\ &\quad + \frac{qf}{(n - f)^2} \|g_0^{(F)'} - g_0^{(E)'}\|_2^2 \\ &= \lambda_{\max}^{\{1, N, E\}} + \frac{q(n - 2f - q)}{(n - f)^2} C^2 + \frac{qf}{(n - f)^2} C^2 \left(\left(\frac{3 + \psi}{2} \right)^2 - \left(\frac{1 + \psi}{2} \right)^2 \right) \\ &> \lambda_{\max}^{\{1, N, E\}}. \end{aligned}$$

(C2) As shown in (C1), the core tension in the choice of the SMEA aggregation rule ultimately reduces to a comparison between the sets $\{1, N, E\}$ and $\{1, N, F\}$, a consideration that also applies for (C2). By Lemma D.12, for all $t \in \{0, \dots, T-1\}$,

$$S_t^* = \{1, N, F\} \iff \left(n - 2f - \frac{2}{m}\right) |g_t^{(F)}|^2 < (n - 2f) |g_t^{(E)}|^2 + \frac{2}{m} |g_t^{(F)}| |g_t^{(E)}|.$$

As $|g_t^{(E)}| < |g_t^{(F)}|$, the above inequality holds if:

$$\left(n - 2f - \frac{2}{m}\right) |g_t^{(F)}|^2 < (n - 2f + \frac{2}{m}) |g_t^{(E)}|^2,$$

which holds by construction as:

$$\psi |g_t^{(F)}| = \psi C < \frac{1 + \psi}{2} C = |g_t^{(E)}|.$$

Hence, by construction SMEA always selects $\{1, N, F\}$ for the unperturbed run, and $\{1, N, E\}$ for the perturbed one. Let denote $p = \frac{f+1/m}{n-f}$, as $\theta_0 = \theta'_0 = 0$, we then have by the (C1) and (C2) conditions,

$$\theta'_T = -\gamma \sum_{t=0}^{T-1} \frac{1}{n-f} \sum_{i \in \{1, N, E\}} g_t^{(i)'} = -\frac{f}{n-f} \gamma \frac{1 + \psi}{2} \gamma C T,$$

and

$$\theta_T = -\gamma \sum_{t=0}^{T-1} \frac{1}{n-f} \sum_{i \in \{1, N, F\}} g_t^{(i)} = p \gamma C T.$$

(iii) Since $\ell(\theta'_T; z) \leq 0$, this directly yields the following lower bound on uniform stability:

$$\sup_{z \in [-C, C]} |\ell(\theta_T; z) - \ell(\theta'_T; z)| \geq \ell(\theta_T; C) = C \theta_T = p \gamma C^2 T.$$

□

D.5.2 Strongly convex case

The structure of the proof with strongly convex objectives closely parallels that of Theorem D.13. We refer to it where appropriate to avoid unnecessary repetition, while explicitly highlighting the key differences.

Theorem D.14 (Added instability in robust GD - strongly convex case). *Consider the setting described in Section 2 under data poisoning attacks. Let $\mathcal{F}_{\mathcal{Z}}$ denote the set of functions that are C -Lipschitz, L -smooth and μ -strongly convex on \mathcal{Z} . Let $\mathcal{A} = \text{GD}$, with the SMEA. Suppose \mathcal{A} is run for $T \in \mathbb{N}^*$ iterations, $T \geq \frac{\ln(1-c)}{\ln(1-\gamma\mu)}$, $0 < c < 1$, with $\gamma \leq \frac{1}{L}$. Then there exists a loss function $\ell \in \mathcal{F}_{\mathcal{Z}}$ and a pair of neighboring datasets such that the uniform stability is lower bounded by*

$$\Omega \left(\frac{C^2}{\mu} \left(\frac{f}{n-f} + \frac{1}{(n-f)m} \right) \right).$$

Proof. In what follows, we provide the proof for the case $d = 1$. The result extends naturally to higher dimensions by embedding the one-dimensional construction into \mathbb{R}^d . We arbitrarily refer to one trajectory as *unperturbed*, with parameter θ_t , and the other as *perturbed*, with θ'_t . We adopt a reasoning analogous to that used in the lower bound proof in Theorem D.13. As explained in Theorem D.13, we structure our proof into three paragraphs (i), (ii) and (iii) below.

(i) Let $C, L, \mu \in \mathbb{R}_+^*$, with $\mu \leq L$, we define the following loss function:

$$\forall \theta \in \Theta = B(0, \frac{C}{2\mu}), z \in \mathcal{Z} = B(0, \frac{C}{2\mu}),$$

$$\ell(\theta; z) = \frac{\mu}{2}(\theta - z)^2, \quad \ell'(\theta; z) = \mu(\theta - z), \quad \ell''(\theta; z) = \mu \quad (19)$$

By construction, the loss function is μ -convex, C -Lipschitz and L -smooth. Let $m, n, f \in \mathbb{N}^*$ such that $f < \frac{n}{2}$. We define,

$$\max\left\{\sqrt{\frac{n - 2(f + \frac{1}{m})}{n - 2(f - \frac{1}{m})}}, 1 - \frac{4}{m(n - 2f)}\right\} < \psi < 1,$$

and the following pair of neighboring datasets:

$$\begin{aligned} S &= \{D_1, \dots, D_n\}, \quad S' = S \setminus \{D_1\} \cup \{D_1 \setminus \{(x^{(1,1)}, y^{(1,1)})\} \cup \{(x'^{(1,1)}, y'^{(1,1)})\}\} \\ \text{with } \forall i \in \{1, \dots, n\}, \quad D_i &= \{(x^{(i,1)}, y^{(i,1)}), \dots, (x^{(i,m)}, y^{(i,m)})\} \\ \text{and } \forall j \in \{1, \dots, m\}, \quad i = 1: z^{(i,j)} &= 0, \quad z^{(1,1)} = \frac{C}{2\mu}, \quad z'^{(1,1)} = 0. \\ i \in N &= \{2, \dots, n - 2f\}, |N| = n - 2f - 1: z^{(i,j)} = 0. \\ i \in E &= \{n - 2f + 1, \dots, n - f\}, |E| = f: z^{(i,j)} = -\frac{1 + \psi}{2} \frac{C}{2\mu}. \\ i \in F &= \{n - f + 1, \dots, n\}, |F| = f: z^{(i,j)} = \frac{C}{2\mu}. \end{aligned} \quad (20)$$

Let $\gamma \leq \frac{1}{L}$ the learning rate. Following this setting, we first state the form that take the parameters and the gradients during the training process, then the conditions under which we claim our instability: for $t \in \{0, \dots, T - 1\}$, $\theta_0 = \theta'_0 = 0$, $g_t^{(N)} = \mu\theta'_t g_t'^{(N)} = g_t'^{(1)} = \mu\theta'_t$, $g_t^{(1)} = (1 - \frac{1}{m})g_t^{(N)} + g_t^{(F)}/m$, $g_t^{(E)} = \mu\left(\theta_t + \frac{1+\psi}{2} \frac{C}{2\mu}\right)$, $g_t^{(E)'} = \mu\left(\theta'_t + \frac{1+\psi}{2} \frac{C}{2\mu}\right)$, $g_t^{(F)} = \mu\left(\theta_t - \frac{C}{2\mu}\right)$ and $g_t^{(F)'} = \mu\left(\theta'_t - \frac{C}{2\mu}\right)$.

(ii) In what follows, we denote, for a subset $S \subset \{1, \dots, n\}$ of size $n - f$. Recall that SMEA performs averaging over the subset of gradients selected according to the following rule, for all $t \in \{0, \dots, T - 1\}$:

$$S_t^* \in \underset{S \subseteq \{1, \dots, n\}, |S| = n - f}{\operatorname{argmin}} \lambda_{\max} \left(\frac{1}{|S|} \sum_{i \in S} (g_t^{(i)} - \bar{g}_S) (g_t^{(i)} - \bar{g}_S)^\top \right).$$

With a slight abuse of notation, we denote, for instance, the set $\{1\} \cup N \cup E$ by $\{1, N, E\}$. To uncover the source of instability, we consider the following two key conditions:

(C1) As shown in the (C1) condition of Theorem D.13, the core tension in the choice of the SMEA aggregation rule ultimately reduces to a comparison between the sets $\{1, N, E\}$ and $\{1, N, F\}$, a consideration that also applies here. By Lemma D.12,

$$S_0^* = \{1, N, E\} \xLeftrightarrow[\text{SMEA}] |g_0^{(F)'}|_2 = \mu \frac{C}{2\mu} > \mu \frac{1 + \psi}{2} \frac{C}{2\mu} = |g_0^{(E)'}|_2.$$

Let denote

$$\lambda_{\max, t}^{\{1, N, E\}'} := \lambda_{\max} \left(\frac{1}{n - f} \sum_{i \in \{1, N, E\}} (g_t^{(i)'} - \bar{g}_S') (g_t^{(i)'} - \bar{g}_S')^\top \right),$$

and

$$\lambda_{\max, t}^{\{1, N, F\}'} := \lambda_{\max} \left(\frac{1}{n - f} \sum_{i \in \{1, N, F\}} (g_t^{(i)'} - \bar{g}_S') (g_t^{(i)'} - \bar{g}_S')^\top \right).$$

We then prove that for all $t \geq 0$, $S_t^{*'} = \{1, N, E\}$, and the unperturbed run converges to the minimizer of the subgroup E . For this property to be true, we require:

$$\lambda_{\max, t}^{\{1, N, E\}'} < \lambda_{\max, t}^{\{1, N, F\}'}.$$

Again, by Lemma D.12, it is equivalent to:

$$\left(\frac{1+\psi}{2}\right)^2 \frac{C^2}{4} < \frac{C^2}{4},$$

which holds by construction. As the condition is independent of t , hence true for all t , this conclude the first condition.

(C2) By Lemma D.12, $S_0^* = \{1, N, F\} \xleftarrow{\text{SMEA and (C1)}} \psi |g_0^{(F)}| = \psi \mu \frac{C}{2\mu} < \mu \frac{1+\psi}{2} \frac{C}{2\mu} = |g_0^{(E)}|$.

We then prove that for all $t \geq 0$, $S_t^* = \{1, N, F\}$, and the perturbed run converges to the minimizer of the subgroup F . For this property to be true, we require:

$$\lambda_{\max, t}^{\{1, N, F\}} < \lambda_{\max, t}^{\{1, N, E\}}.$$

Again, by Lemma D.12, it is equivalent to:

$$\begin{aligned} & f(n-2f-1) \frac{C^2}{4} + f \frac{C^2}{4} \left(1 - \frac{1}{m}\right)^2 + (n-2f-1) \frac{C^2}{4m^2} \\ & < f(n-2f-1) \left(\frac{1+\psi}{2}\right)^2 \frac{C^2}{4} + f \frac{C^2}{4} \left(\frac{1+\psi}{2} + \frac{1}{m}\right)^2 + (n-2f-1) \frac{C^2}{4m^2} \end{aligned}$$

That is,

$$0 < -(n-2f-1) \left(\frac{1-\psi}{2}\right) \left(\frac{3+\psi}{2}\right) + \left(\frac{3+\psi}{2}\right) \left(\frac{2}{m} - \frac{1-\psi}{2}\right),$$

which is equivalent to:

$$\psi > \frac{4}{(n-f)m}.$$

As the condition is independent of t , hence true for all t , this conclude the second condition.

Let denote $p = \frac{f+1/m}{n-f}$. We have proven that we have for $t \geq 0$:

$$\text{SMEA}(g_t^{(1)}, \dots, g_t^{(n)}) = \frac{1}{n-f} \sum_{i \in \{1, N, F\}} g_t^{(i)} = \mu \theta_t - p \frac{C}{2}$$

$$\text{SMEA}(g_t^{(1)'}, \dots, g_t^{(n)'}) = \frac{1}{n-f} \sum_{i \in \{1, N, E\}} g_t^{(i)'} = \mu \theta'_t + \frac{f}{n-f} \frac{1+\psi}{2} \frac{C}{2}$$

For $t \geq 0$, $\theta_{t+1} = (1-\gamma\mu)\theta_t + \frac{\gamma p C}{2}$ and $\theta'_{t+1} = (1-\gamma)\theta'_t - \frac{\gamma f(1+\psi)C}{4(n-f)}$. Hence: $\theta'_T \leq 0$ and,

$$\theta_T = \frac{pC}{2\mu} \left(1 - (1-\gamma\mu)^T\right)$$

(iii) Provided that there exist $0 < c < 1$ such that $T \geq \frac{\ln(1-c)}{\ln(1-\gamma\mu)}$, this implies: $|\theta_T - \theta'_T| \geq |\theta_T| \in \Omega(\frac{pC}{\mu})$. Finally, let define $\chi = -\text{sign}(\theta_T + \theta'_T)$, we have:

$$\begin{aligned} & |\ell(\theta_T; \chi \frac{C}{2\mu}) - \ell(\theta'_T; \chi \frac{C}{2\mu})| = \frac{\mu}{2} \left| \left(\theta_T - \chi \frac{C}{2\mu}\right)^2 - \left(\theta'_T - \chi \frac{C}{2\mu}\right)^2 \right| \\ & = \frac{\mu}{2} |\theta_T - \theta'_T| |\theta_T + \theta'_T + \text{sign}(\theta_T + \theta'_T) \frac{C}{\mu}| \\ & \geq \frac{\mu}{2} |\theta_T| \frac{C}{\mu} \in \Omega\left(\frac{pC^2}{\mu}\right) \end{aligned}$$

□

D.6 Proof of Theorem 3.3 for projected SGD

In this section, we derive a lower bound on the uniform stability of distributed projected SGD with the SMEA aggregation rule under poisoning attacks. The structure of the proof closely parallels that of Theorem D.13. We refer to it where appropriate to avoid unnecessary repetition, while explicitly highlighting the key differences.

Theorem D.15 (Added instability in robust projected SGD). *Consider the setting described in Section 2 under data poisoning attacks. Let $\mathcal{F}_{\mathcal{Z}}$ denote the set of functions that are C -Lipschitz, L -smooth and convex on \mathcal{Z} . Let $\mathcal{A} = \text{projected-SGD}$, with the SMEA. Suppose \mathcal{A} is run for $T \in \mathbb{N}^*$ iterations with $\gamma \leq \frac{1}{L}$. Assume there exist a constant $\tau \geq c > 0$ for an arbitrary c , such that $T \geq \tau m$, then there exists a loss function $\ell \in \mathcal{F}_{\mathcal{Z}}$ and a pair of neighboring datasets such that the uniform stability is lower bounded by*

$$\Omega \left(\gamma C^2 T \left(\frac{f}{n-f} + \frac{1}{(n-f)m} \right) \right).$$

Proof. We adopt a reasoning analogous to that used in the lower bound proof in Theorem D.13. However, we can no longer rely on a linear loss function. The reason is that, in the linear case, the gradients are independent of the model parameters, so when the differing samples are not selected, the perturbed and unperturbed executions remain indistinguishable—thus precluding any divergence. We also have to account for the first occurrence of the differing samples, denoted as $T_0 = \inf\{t; J_{t-1}^{(1)} = 1\}$. In fact, we require to ensure the conditions (C1) (that is $S_0^* = \{1, N, E\}$) at initialization but we now would like to ensure (C2) when the differing samples are drawn (that is $S_t^* = \{1, N, F\}$ for $t \geq t_0 - 1$, knowing $T_0 = t_0$ for any time $t_0 < T$). To do so, we modify the possible parameters we are searching through our optimization process, making the condition at initialization sufficient to ensure (C1) and (C2) for any $t_0 \in \{0, \dots, T-1\}$. By defining $\Theta = \{\lambda v; 0 \leq \lambda < \infty\}$ and applying the projection p_{Θ} at each step to constrain the parameter within the set of admissible values, we ensure, under condition (C1) that $\theta'_t = 0$ for all $t > 0$ and $\theta_t = 0$ for $t < t_0$. Consequently, the condition (C2) is ensured by construction for any time the differing sample is drawn. While using a projection simplifies the lower bound proof (Appendix D.6), it does not affect the tightness of our result since the upper bound on uniform stability remains valid for projected-SGD as the projection does not increase the distance between projected points.

As explained in Theorem D.13, we structure our proof into three paragraphs (i), (ii) and (iii) below.

(i) Let $C, L \in \mathbb{R}_+^*$, we define the following dataset and loss:

$$\begin{aligned} \forall \theta \in \Theta = \{\lambda v; \lambda \in \mathbb{R}\} \subset \mathbb{R}^2, (x, y) \in B(0, \sqrt{L}) \times \mathbb{R}, \\ \ell(\theta; (x, y)) &= \begin{cases} \frac{1}{2}(\theta^T x - y)^2 & \text{if } |\theta^T x - y| \|x\|_2 \leq C \\ C \left(\frac{|\theta^T x - y|}{\|x\|_2} - \frac{C}{2\|x\|_2^2} \right) & \text{otherwise.} \end{cases} \\ \nabla \ell(\theta; (x, y)) &= \begin{cases} (\theta^T x - y)x & \text{if } |\theta^T x - y| \|x\|_2 \leq C \\ C \operatorname{sign}(\theta^T x - y) \frac{x}{\|x\|_2} & \text{otherwise.} \end{cases} \\ \nabla^2 \ell(\theta; (x, y)) &= \begin{cases} 0 \preceq xx^T \preceq LI_d & \text{if } |\theta^T x - y| \|x\|_2 \leq C \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

By construction, the loss function is convex, C -Lipschitz and L -smooth. Let $m, n, f \in \mathbb{N}^*$ such that $f < \frac{n}{2}$. We define,

$$\begin{aligned} 0 < \psi &= \frac{n-2f-2}{n-2f} < 1, \\ \min\{1 - \psi, \gamma L \frac{f}{n-f}\} &> \epsilon > 0 \end{aligned}$$

and the following pair of neighboring datasets:

$$\begin{aligned}
v &\in B(0, \sqrt{L}), \quad \|v\|_2^2 = L, \quad \beta = \frac{C}{\sqrt{L}}, \quad \alpha = (1 - \epsilon)\beta, \quad b = \frac{1}{\sqrt{T}} \\
S &= \{D_1, \dots, D_n\}, \quad S' = S \setminus \{D_1\} \cup \{D_1 \setminus \{(x^{(1,1)}, y^{(1,1)})\} \cup \{(x'^{(1,1)}, y'^{(1,1)})\}\} \\
\text{with } \forall i \in \{1, \dots, n\}, \quad D_i &= \{(x^{(i,1)}, y^{(i,1)}), \dots, (x^{(i,m)}, y^{(i,m)})\} \\
\text{and } \forall j \in \{1, \dots, m\}, \quad i = 1: &(x^{(1,j)}, y^{(1,j)}) = (0, 0), \quad (x^{(1,1)}, y^{(1,1)}) = (bv, \frac{\beta}{b}), \\
&\text{and } (x'^{(1,1)}, y'^{(1,1)}) = (0, 0). \\
i \in N = \{2, \dots, n - 2f\}, |N| &= n - 2f - 1: (x^{(i,j)}, y^{(i,j)}) = (0, 0). \\
i \in E = \{n - 2f + 1, \dots, n - f\}, |E| &= f: (x^{(i,j)}, y^{(i,j)}) = (v, -\alpha). \\
i \in F = \{n - f + 1, \dots, n\}, |F| &= f: (x^{(i,j)}, y^{(i,j)}) = (bv, \frac{\beta}{b}).
\end{aligned}$$

Let $\gamma \leq \frac{1}{L}$ the learning rate. Following this setup, we first describe the form of the parameters and gradients throughout the training process, and then state the conditions under which we establish our instability result, assuming $\theta_0 = \theta'_0 = 0$. For $t \in \{0, \dots, T - 1\}$, and $\theta_t = \lambda_t v$, $\theta'_t = \lambda'_t v$, $\lambda_t, \lambda'_t \in [\lambda_E^* := -\frac{\alpha}{L}, \lambda_F^* := \frac{\beta}{b^2 L}]$:

$$\begin{aligned}
g_t^{(N)} &= 0, \quad \|g_t^{(E)}\|_2 > \psi \|g_t^{(F)}\|_2, \quad g_t^{(F)} = - \underbrace{(\beta - \lambda_t L b^2)}_{>0 \text{ and } < \frac{C}{\sqrt{L}}} v \\
g_t^{(N)'} &= 0, \quad g_t^{(E)'} = \underbrace{(\lambda'_t L + \alpha)}_{>0 \text{ and } < \frac{C}{\sqrt{L}}} v, \quad \|g_t^{(F)'}\|_2 > \|g_t^{(E)'}\|_2
\end{aligned}$$

(ii) In what follows, we denote, for a subset $S \subset \{1, \dots, n\}$ of size $n - f$. Recall that SMEA performs averaging over the subset of gradients selected according to the following rule, for all $t \in \{0, \dots, T - 1\}$:

$$S_t^* \in \underset{S \subseteq \{1, \dots, n\}, |S|=n-f}{\operatorname{argmin}} \lambda_{\max} \left(\frac{1}{|S|} \sum_{i \in S} (g_t^{(i)} - \bar{g}_S) (g_t^{(i)} - \bar{g}_S)^\top \right).$$

With a slight abuse of notation, we denote, for instance, the set $\{1\} \cup N \cup E$ by $\{1, N, E\}$. To uncover the source of instability, we consider the following two key conditions:

(C1) As shown in the (C1) condition of Theorem D.13, the core tension in the choice of the SMEA aggregation rule ultimately reduces to a comparison between the sets $\{1, N, E\}$ and $\{1, N, F\}$, a consideration that also applies here. $S_0^{*'} = \{1, N, E\} \xLeftrightarrow{\text{SMEA}} \|g_0^{(F)'}\|_2 = \beta\sqrt{L} > \alpha\sqrt{L} = \|g_0^{(E)'}\|_2$. Hence, $\theta'_1 = p_\Theta(-\gamma \frac{f}{n-f} g_0^{(E)}) = p_\Theta(-\gamma \frac{f}{n-f} \alpha v) = 0$. We then prove by induction, that for all $t \geq 0$, $S_t^{*'} = \{1, N, E\}$ and the unperturbed run remains at 0. The perturbed run remains at 0 until the differing samples are drawn.

(C2) Let denote

$$\lambda_{\max, t}^{\{1, N, E\}} := \lambda_{\max} \left(\frac{1}{n-f} \sum_{i \in \{1, N, E\}} (g_t^{(i)} - \bar{g}_S) (g_t^{(i)} - \bar{g}_S)^\top \right),$$

and

$$\lambda_{\max, t}^{\{1, N, F\}} := \lambda_{\max} \left(\frac{1}{n-f} \sum_{i \in \{1, N, F\}} (g_t^{(i)} - \bar{g}_S) (g_t^{(i)} - \bar{g}_S)^\top \right).$$

Assuming $T_0 = t_0$ known, $S_{t_0-1}^* = \{1, N, F\} \xLeftrightarrow[\text{SMEA and (C1)}]{\frac{n-2(f+1)}{n-2}} \|g_{t_0-1}^{(F)}\|_2 < \|g_{t_0-1}^{(E)}\|_2$. In fact, by Lemma D.12 with $g_A = g_{t_0-1}^{(E)}$, $g_B = g_{t_0-1}^{(F)}$ and $g_C = 0$:

$$S_{t_0-1}^* \in \operatorname{argmin} \left\{ \lambda_{\max, t_0-1}^{\{1, N, F\}} = \frac{(f+1)(n-2f-1)}{(n-f)^2} \|g_{t_0-1}^{(F)}\|_2^2, \right. \\ \lambda_{\max, t_0-1}^{\{1, N, E\}} = \frac{f(n-2f-1)}{(n-f)^2} \|g_{t_0-1}^{(E)}\|_2^2 + \frac{(n-2f-1)}{(n-f)^2} \|g_{t_0-1}^{(F)}\|_2^2 \\ \left. + \frac{f}{(n-f)^2} \left(\|g_{t_0-1}^{(E)}\|_2 + \|g_{t_0-1}^{(F)}\|_2 \right)^2 \right\}.$$

Hence,

$$S_{t_0-1}^* = \{1, N, F\} \Leftrightarrow (n-2f-1) (\|g_{t_0-1}^{(F)}\|_2^2 - \|g_{t_0-1}^{(E)}\|_2^2) < (\|g_{t_0-1}^{(F)}\|_2 + \|g_{t_0-1}^{(E)}\|_2)^2 \\ \Leftrightarrow (n-2f-1) (\|g_{t_0-1}^{(F)}\|_2 - \|g_{t_0-1}^{(E)}\|_2) < \|g_{t_0-1}^{(F)}\|_2 + \|g_{t_0-1}^{(E)}\|_2 \\ \Leftrightarrow \psi \|g_{t_0-1}^{(F)}\|_2 < \|g_{t_0-1}^{(E)}\|_2$$

which holds by construction of our initialization, $\psi \|g_{t_0-1}^{(F)}\|_2 < \|g_{t_0-1}^{(E)}\|_2 < \|g_{t_0-1}^{(F)}\|_2$.

We then prove that for the next step, $S_t^* = \{1, N, F\}$ in any case. This is sufficient to prove by induction that for all $t \geq t_0$, $S_t^* = \{1, N, F\}$ and the perturbed run converges to the minimizer of the subgroup F . In fact, for $t \geq t_0$, $\|g_t^{(F)}\|_2$ decreases and $\|g_t^{(E)}\|_2$ increases.

- (1) with probability $\frac{1}{m}$, the differing samples are again drawn. As we already ensured the conditions in the previous step, we have $\psi \|g_{t_0}^{(F)}\|_2 \leq \|g_{t_0-1}^{(F)}\|_2 < \|g_{t_0-1}^{(E)}\|_2 \leq \|g_{t_0}^{(E)}\|_2$, then $S_{t_0}^* = \{1, N, F\}$.
- (2) with probability $1 - \frac{1}{m}$ we have to ensure $S_{t_0}^* = \{1, N, F\} \xLeftrightarrow{\text{SMEA}} \|g_{t_0}^{(F)}\|_2 \leq \|g_{t_0}^{(E)}\|_2$ with $\theta_{t_0} = -\gamma \frac{f+1}{n-f} g_{t_0-1}^{(F)} = \gamma \beta \frac{f+1}{n-f} v$ hence:

$$\|g_{t_0}^{(F)}\|_2 = -\beta \left(1 - \gamma \frac{f+1}{n-f} Lb^2 \right) v, \\ \|g_{t_0}^{(E)}\|_2 = \max\{\beta, \beta \left(1 - \epsilon + \gamma L \frac{f+1}{n-f} \right)\} v = \beta v$$

We proved that for $t \geq t_0 - 1$, $S_{t_0-1}^* = \{1, N, F\}$. Hence:

$$\mathbb{E}[\theta_t | T_0 = t_0] = \begin{cases} 0 & \text{if } t \leq t_0 - 1 \\ \gamma \beta \frac{f+1}{n-f} v & \text{for } t = t_0 \\ \mathbb{E}[\theta_{t-1} | T_0 = t_0] - \gamma p \mathbb{E}[g_{t-1}^{(F)} | T_0 = t_0] & \text{otherwise.} \end{cases}$$

with $\theta_t = \lambda_t v$, $p = \frac{f+1}{n-f}$, the recurrence for $t > t_0$ is given by:

$$\mathbb{E}[\lambda_t | T_0 = t_0] = \mathbb{E}[\lambda_{t-1} | T_0 = t_0] + p\gamma(\beta - \mathbb{E}[\lambda_{t-1} | T_0 = t_0] Lb^2) \\ = (1 - p\gamma Lb^2) \mathbb{E}[\lambda_{t-1} | T_0 = t_0] + p\gamma\beta$$

Hence:

$$\mathbb{E}[\theta_t | T_0 = t_0] = \begin{cases} 0 & \text{if } t \leq t_0 - 1 \\ \frac{f+1}{f+\frac{1}{m}} \frac{\beta}{b^2 L} \left(1 - (1 - p\gamma Lb^2)^{t-t_0+1} \right) v & \text{otherwise.} \end{cases}$$

Taking the total expectation, we obtain:

$$\mathbb{E}[\theta_T] = \sum_{t_0=1}^T \mathbb{E}[\theta_T | T_0 = t_0] \mathbb{P}(T_0 = t_0) \\ = \frac{CT \frac{f+1}{f+\frac{1}{m}}}{Lm} \sum_{t_0=1}^T \left(1 - \left(1 - \frac{p\gamma L}{T} \right)^{T-t_0+1} \right) \left(1 - \frac{1}{m} \right)^{t_0-1} \frac{v}{\sqrt{L}}$$

(iii) In the remainder of the proof, we bound $\mathbb{E}[\theta_T]$. Specifically, we leverage its analytical expression to establish a concise lower bound on uniform argument stability.

$$\begin{aligned}
& \sum_{t_0=1}^T \left(1 - (1 - p\gamma Lb^2)^{T-t_0+1}\right) \left(1 - \frac{1}{m}\right)^{t_0-1} \\
& \geq \sum_{t_0=1}^T \left(1 - \frac{1}{1 + p\gamma Lb^2(T-t_0+1)}\right) \left(1 - \frac{1}{m}\right)^{t_0-1} \\
& \geq \sum_{t_0=1}^T \frac{p\gamma L \frac{T-t_0+1}{T}}{1 + p\gamma L \frac{T-t_0+1}{T}} \left(1 - \frac{1}{m}\right)^{t_0-1}
\end{aligned}$$

Where we used that $(1 - p\gamma Lb^2)^{T-t_0+1} \leq \frac{1}{1 + p\gamma Lb^2(T-t_0+1)}$ for $T - t_0 + 1 \geq 0$ and $-1 < -p\gamma Lb^2 < 0$. In fact using $\ln(1+x) \leq x$ for $x > -1$,

$$\begin{aligned}
\ln \left[(1 - p\gamma Lb^2)^{T-t_0+1} (1 + p\gamma Lb^2(T-t_0+1)) \right] &= (T-t_0+1) \ln(1 - p\gamma Lb^2) \\
&\quad + \ln(1 + p\gamma Lb^2(T-t_0+1)) \\
&\leq -p\gamma Lb^2(T-t_0+1) + p\gamma Lb^2(T-t_0+1) \\
&= 0.
\end{aligned}$$

Also, as $\frac{1}{1 + p\gamma L \frac{T-t_0+1}{T}} \geq \frac{1}{1 + \gamma Lp} \geq \frac{1}{1+p} \geq \frac{1}{2}$, we have:

$$\begin{aligned}
& \sum_{t_0=1}^T \left(1 - (1 - p\gamma Lb^2)^{T-t_0+1}\right) \left(1 - \frac{1}{m}\right)^{t_0-1} \\
& \geq \frac{p\gamma L}{2} \sum_{t_0=1}^T \left(1 - \frac{t_0-1}{T}\right) \left(1 - \frac{1}{m}\right)^{t_0-1} \\
& = \frac{p\gamma Lm}{2} \left(1 - \left(1 - \frac{1}{m}\right)^T\right) - \frac{p\gamma L}{2T} \sum_{t_0=1}^T (t_0-1) \left(1 - \frac{1}{m}\right)^{t_0-1} \\
& = \frac{p\gamma Lm}{2} \left(1 - \left(1 - \frac{1}{m}\right)^T\right) - \frac{p\gamma L}{2T} \left[m(m-1) - m^2 \left(1 - \frac{1}{m}\right)^T \left(1 + \frac{T-1}{m}\right) \right] \\
& = \frac{p\gamma Lm}{2} \left[1 - \frac{m-1}{T} \left(1 - \left(1 - \frac{1}{m}\right)^T\right) \right]
\end{aligned}$$

where we used the finite polylogarithmic sum formula in the second equality above.

If we additionally assume there exist a constant $\tau \geq c > 0$ for an arbitrary c , such that $T \geq \tau m$, we then can prove that $K(T, m) = \left[1 - \frac{m-1}{T} \left(1 - \left(1 - \frac{1}{m}\right)^T\right)\right] \in \Omega(1)$. In fact this is trivial for $m = 1$, then we prove that the quantity $K(T, m)$ increases with T , for $T \geq 1$ and $m \geq 2$. We prove by induction that $K(T+1) - K(T) = \frac{m-1}{T(T+1)} \left[1 - \left(1 - \frac{1}{m}\right)^T \left(1 + \frac{T}{m}\right)\right] > 0$ which is equivalent to $\left(1 - \frac{1}{m}\right)^T \left(1 + \frac{T}{m}\right) < 1$. For $T = 1$, $\left(1 - \frac{1}{m}\right)\left(1 + \frac{1}{m}\right) = 1 - \frac{1}{m^2} < 1$. For $T+1 > 2$, $\frac{\left(1 - \frac{1}{m}\right)^{T+1} \left(1 + \frac{T+1}{m}\right)}{\left(1 - \frac{1}{m}\right)^T \left(1 + \frac{T}{m}\right)} = \left(1 - \frac{1}{m}\right)^{1 + \frac{T+1}{m}} < 1$.

Hence:

$$\begin{aligned}
1 - \frac{m-1}{T} \left(1 - \left(1 - \frac{1}{m}\right)^T\right) &\geq 1 - \frac{m-1}{\tau m} \left(1 - \left(1 - \frac{1}{m}\right)^{\tau m}\right) \\
&\stackrel{(1)}{\geq} 1 - \frac{1}{\tau} \left(1 - \frac{1}{m}\right) \left(1 - e^{-\tau} \left(1 - \frac{\tau}{m}\right)\right)
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(2)}{\geq} 1 - \frac{1 - e^{-\tau}}{\tau} \\
&\geq 1 - \frac{1 - e^{-c}}{c} > 0
\end{aligned}$$

where we used in the second inequality (1) that $(1 - \frac{1}{m})^{\tau m} \geq e^{-\tau}(1 - \frac{\tau}{m})$ (we can prove it by taking the logarithm and compare their Taylor expansion). We also used in the third inequality (2) that $1 - \frac{1}{\tau}(1 - \frac{1}{m})(1 - e^{-\tau}(1 - \frac{\tau}{m})) - 1 - \frac{1 - e^{-\tau}}{\tau} = \frac{1}{\tau m} \left(1 - (1 + \tau)e^{-\tau} + \frac{\tau e^{-\tau}}{m}\right) \geq \frac{e^{-\tau}}{m^2} \geq 0$.

Wrapping-up:

$$\mathbb{E}[\|\theta_T - \theta'_T\|_2] \geq \|\mathbb{E}[\theta_T]\|_2 \geq \frac{1}{2} \left(1 - \frac{1 - e^{-c}}{c}\right) \underbrace{\frac{f+1}{f + \frac{1}{m}}}_{\geq 1} p\gamma CT \in \Omega(p\gamma CT)$$

Building on this lower bound for uniform argument stability, we then establish a connection to uniform stability, allowing us to complete the comparison with the upper bound from Theorem D.11.

$$\mathbb{E} \left[\ell(\theta_T; (v, -\frac{C}{\sqrt{L}})) - \ell(\theta'_T; (v, -\frac{C}{\sqrt{L}})) \right] = C\sqrt{L}\mathbb{E}\lambda_T + \frac{C^2}{2L} - \frac{C^2}{2L} \in \Omega(pC^2\gamma T)$$

□

E Deferred proofs and additional details from Section 4

We provide details supporting Section 4, including numerical experiments E.1 and generalization error computations E.2.

E.1 Numerical analysis additional details

We specify here the details of our numerical experiment presented in Section 4 into three points: (i) the problem instantiation; (ii) our tailored *Byzantine failure* attack; and (iii) how we compute the empirical robustness parameter.

- (i) We instantiate the setting from the proof of theorem D.13 (see full details therein). That is, we consider $\mathcal{A} = \text{GD}$ under a linear loss function, using the parameter $C = 1, \gamma = 1, T = 5, n = 15, m = 1$. We make vary f from 1 to 7, and for each value instantiate the dataset presented in the proof of theorem D.13.
- (ii) In our tailored Byzantine failure attack, the identities of the Byzantine workers remain fixed throughout the optimization process, but vary depending on the number of Byzantine workers f . The set of Byzantine identities is defined as $\mathcal{B} := \{1, \dots, n\} \setminus \mathcal{H}$, and is summarized in the table below. As we will see, we optimize the values sent by Byzantine workers so that they are selected while also biasing the gradient statistics—either toward positive or negative values—depending on the direction of the first parameter. To construct a maximally impactful scenario, we choose the identities of the Byzantine workers to maximize the minimal variance over all subsets of size $n - f$. This increases their biasing power, as it hinges on the smallest variance among these subsets. For instance, when $f = 2$, the Byzantine workers can be hidden in the dominant subset N ; however, when $f = 6$, they must be distributed across both competing subgroups E and F to maintain influence.

$$\begin{array}{ll}
f = 1 : \mathcal{B} = \{2\} & f = 4 : \mathcal{B} = \{2, 3, 4, 5\} \\
f = 2 : \mathcal{B} = \{2, 3\} & f = 5 : \mathcal{B} = \{2, 3, 4, 5, 6\} \\
f = 3 : \mathcal{B} = \{2, 3, 4\} & f = 6 : \mathcal{B} = \{7, 8, 9, 10, 11, 12\} \\
& f = 7 : \mathcal{B} = \{6, 7, 8, 9, 10, 11, 12\}
\end{array}$$

The Byzantine strategy mimics the poisoning attack in the first iteration, then amplifies the resulting parameter direction by adaptively crafting updates that remain within the aggregation rule’s selection range. Specifically, we describe their algorithm in what follows:

- If $\theta = 0$, all Byzantine workers send the same value as the one in the poisoning attack. This value depends on each worker’s identity; see the proof of Theorem D.13 for details.
- If $\theta > 0$, all Byzantine workers send the value $-\alpha_f$, where α_f is numerically chosen to be as large as possible under the constraint that the SMEA aggregation rule continues to select the values sent by Byzantine workers in subsequent steps.
- If $\theta < 0$, all Byzantine workers send the value β_f , where β_f is numerically chosen to be as large as possible under the constraint that the SMEA aggregation rule continues to select the values sent by Byzantine workers in subsequent steps.

The following pseudocode summarizes the behavior of the Byzantine workers. Let $\alpha \in \mathbb{R}$ and

$$S_\alpha^* \in \underset{|S|=n-f}{\operatorname{argmin}} \operatorname{Var} \left(\{g_k\}_{k \in S \cap \mathcal{H}} \cup \{\alpha\}_{k \in S \cap \mathcal{B}} \right).$$

Algorithm 1: Byzantine worker $i \notin \mathcal{H}$ behavior, $\epsilon = 10^{-3}$

```

if  $\theta = 0$  then
    return  $g_{\text{poisoning}}^{(i)}$  # the same value as data poisoning
else if  $\theta > 0$  then
    return  $\alpha_f = \inf \left\{ \alpha; i \in S_\alpha^* \right\} + \epsilon$  # amplifies the initial direction
else if  $\theta < 0$  then
    return  $\beta_f = \sup \left\{ \beta; i \in S_\beta^* \right\} - \epsilon$  # amplifies the initial direction

```

- (iii) Let denote $\hat{\kappa}_f, \hat{\kappa}'_f$ the empirical measures of the robustness coefficients observed over the optimization trajectories. The empirical robustness coefficients are computed as follows:

$$\hat{\kappa}_f = \max_{t \in \{0, \dots, T-1\}} \max_{\substack{S \subset \{1, \dots, n\} \\ |S|=n-f}} \left(\frac{|\bar{g}_{S_t^*} - \bar{g}_S|^2}{\frac{1}{|S|} \sum_{i \in S} \left(g_t^{(i)} - \bar{g}_S \right)^2} \right),$$

and

$$\hat{\kappa}'_f = \max_{t \in \{0, \dots, T-1\}} \max_{\substack{S \subset \{1, \dots, n\} \\ |S|=n-f}} \left(\frac{|\bar{g}'_{S_t'^*} - \bar{g}'_S|^2}{\frac{1}{|S|} \sum_{i \in S} \left(g_t'^{(i)} - \bar{g}'_S \right)^2} \right),$$

where $\bar{g}_S = \frac{1}{|S|} \sum_{i \in S} g_t^{(i)}$ and S_t^* the subset chosen by the SMEA aggregation rule at iteration t . The same notation applies to the other run. To assess the gap between our theoretical upper bound and its empirical realization, we analyze the tightness of the robustness coefficient $\hat{\kappa}_f$ as observed in our setup. We find that $\hat{\kappa}_f$ is consistently smaller than its theoretical worst-case value $\kappa_{\text{SMEA}} = \frac{4f}{n-f} \left(1 + \frac{f}{n-2f} \right)^2$. Substituting this empirical value into our theoretical upper bound, $\gamma CT \left(\frac{2}{(n-f)m} + 2\sqrt{\kappa} \right)$, significantly reduces the discrepancy between theoretical predictions and actual performance:

$$\gamma CT \left(\frac{2}{(n-f)m} + \sqrt{\hat{\kappa}_f} + \sqrt{\hat{\kappa}'_f} \right).$$

E.2 Details on the computation of the generalization error

We again instantiate the problem outlined in the proof of our lower bound (Theorem 3.3, see Appendix D.5.1), under data poisoning attack, with $\gamma = C = 1$ and $m = 1$. To calculate the generalization error, we assume every honest local distribution is a Dirac with singularity at the value defined in the proof of Theorem 3.3—except for a “pivot” worker (indexed 1 without loss of generality) whose distribution is a mixture $p_1 = \frac{1}{2}\delta_0 + \frac{1}{2}\delta_{-1}$. In this case, the two possible datasets are always neighboring. We denote by $\theta_T^{(0)}$ and $\theta_T^{(-1)}$ the parameters output by the robust distributed

GD algorithm with SMEA on the two different datasets. We denote by \mathcal{S} a dataset sampled from the distribution specified above, and let z_1 be the sample held by worker 1. In this setup, we are able to compute the generalization error exactly as follows:

$$\begin{aligned}
\mathbb{E}_{\mathcal{S}}[R_{\mathcal{H}}(\mathcal{A}(\mathcal{S})) - \widehat{R}_{\mathcal{H}}(\mathcal{A}(\mathcal{S}))] &= \frac{1}{|\mathcal{H}|} \mathbb{E}_{\mathcal{S}}[\mathbb{E}_{z \sim p_1}[\theta_T^{(z_1)} z] - \theta_T^{(z_1)} z_1] \\
&= \frac{1}{|\mathcal{H}|} \mathbb{E}_{\mathcal{S}}[\frac{1}{2} \theta_T^{(z_1)} \times 0 + \frac{1}{2} \theta_T^{(z_1)} \times (-1) - \theta_T^{(z_1)} z_1] \\
&= \frac{1}{|\mathcal{H}|} \left(-\frac{1}{2} \mathbb{E}_{\mathcal{S}} \theta_T^{(z_1)} - \mathbb{E}_{\mathcal{S}}[\theta_T^{(z_1)} z_1] \right) \\
&= \frac{1}{|\mathcal{H}|} \left(-\frac{1}{2} \left(\frac{1}{2} \theta_T^{(0)} + \frac{1}{2} \theta_T^{(-1)} \right) - \frac{1}{2} \theta_T^{(0)} \times 0 - \frac{1}{2} \theta_T^{(-1)} \times (-1) \right) \\
&= \frac{1}{|\mathcal{H}|} \left(-\frac{1}{4} \theta_T^{(0)} - \frac{1}{4} \theta_T^{(-1)} + \frac{1}{2} \theta_T^{(-1)} \right) \\
&= \frac{\theta_T^{(-1)} - \theta_T^{(0)}}{4|\mathcal{H}|}
\end{aligned}$$

Hence:

$$\mathbb{E}_{\mathcal{S}}[R_{\mathcal{H}}(\mathcal{A}(\mathcal{S})) - \widehat{R}_{\mathcal{H}}(\mathcal{A}(\mathcal{S}))] = \frac{\theta_T^{(-1)} - \theta_T^{(0)}}{4(n-f)}.$$

This value corresponds to the stability plots shown in Figure 1, modulo the factor $\frac{1}{4(n-f)}$. Hence, in the worst-case, the generalization error under Byzantine failure and data poisoning attacks exhibits the same discrepancy as the one observed in uniform algorithmic stability.

F Refined bounds under bounded heterogeneity and bounded variance assumptions

For clarity and simplicity, we chose to emphasize our results under the broad bounded gradient assumption (*i.e.* Lipschitz-continuity of the loss function) in the main text. In this section, we present results derived under refined assumptions—namely, bounded heterogeneity and bounded variance—in place of the more general bounded gradient condition. These assumptions are more commonly used in the optimization literature (as opposed to the generalization literature).

We formally define these assumptions below.

Assumption F.1 (Bounded heterogeneity). *There exists $G < \infty$ such that for all $\theta \in \Theta$: $\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \|\nabla \widehat{R}_i(\theta) - \nabla \widehat{R}_{\mathcal{H}}(\theta)\|_2^2 \leq G^2$.*

Assumption F.2 (Bounded variance). *There exists $\sigma < \infty$ such that for each honest workers $i \in \mathcal{H}$ and all $\theta \in \Theta$: $\frac{1}{m} \sum_{x \in \mathcal{D}_i} \|\nabla_{\theta} \ell(\theta; x) - \nabla \widehat{R}_i(\theta)\|_2^2 \leq \sigma^2$.*

As shown in the following, our analysis techniques can yield tighter bounds when assuming bounded heterogeneity and bounded variance. However, we note that these bounds do not provide additional conceptual insights within the scope of our discussion. Indeed, within the uniform stability framework, if we assume only bounded gradients and nothing more, there exist distributed learning settings where G is a constant multiple of the Lipschitz constant C . For instance, consider a scenario where half of the honest workers produce gradients with norm C , and the other half produce zero gradients; in this case, G would be $C/2$.

As an example, we derive bounds on the expected trace and spectral norm of the empirical covariance matrix of honest workers' gradients that are expressed in terms of G and σ . These estimates can be directly used in Theorems D.3 and D.5 to provide tighter bounds under additional assumptions of bounded heterogeneity and bounded variance.

Lemma F.1 (Expectation bound for the empirical covariance trace). *Consider the setting described in Section 2 under Byzantine attacks with Assumptions F.1 and F.2. Let $\mathcal{A} \in \{\text{GD}, \text{SGD}\}$, with a (f, κ, Tr) -robust aggregation rule, for $T \in \mathbb{N}^*$ iterations. We have for every $t \in \{1, \dots, T-1\}$:*

$$\mathbb{E}_{\mathcal{A}} \text{Tr} \Sigma_{\mathcal{H}, t} \leq 3G^2 + 3\sigma^2 \left(1 + \frac{1}{n-f}\right)$$

Proof. We decompose the sum into the following:

$$\begin{aligned}
\mathbb{E}_{\mathcal{A}} \text{Tr } \Sigma_{\mathcal{H},t} &= \mathbb{E} \left[\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \|g_t^{(i)} - \nabla \widehat{R}_i(\theta_t) + \nabla \widehat{R}_i(\theta_t) - \nabla \widehat{R}_{\mathcal{H}}(\theta_t) + \nabla \widehat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t^{(i)}\|_2^2 \right] \\
&\leq \frac{3}{n-f} \sum_{i \in \mathcal{H}} \mathbb{E}[\|g_t^{(i)} - \nabla \widehat{R}_i(\theta_t)\|_2^2] + 3\mathbb{E}[\frac{1}{n-f} \sum_{i \in \mathcal{H}} \|\nabla \widehat{R}_i(\theta_t) - \nabla \widehat{R}_{\mathcal{H}}(\theta_t)\|_2^2] \\
&\quad + \frac{3}{n-f} \sum_{i \in \mathcal{H}} \mathbb{E}[\|\nabla \widehat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t^{(i)}\|_2^2] \\
&\leq 3G^2 + 3\sigma^2 + \frac{3}{(n-f)^2} \sum_{i \in \mathcal{H}} \mathbb{E}[\|g_t^{(i)} - \nabla \widehat{R}_i(\theta_t)\|_2^2]
\end{aligned}$$

□

Lemma F.2 (Expectation bound for the empirical covariance spectral norm). *Consider the setting described in Section 2 under Byzantine attacks with Assumptions F.1 and F.2. Let $\mathcal{A} \in \{\text{GD}, \text{SGD}\}$, with a $(f, \kappa, \|\cdot\|_{sp})$ -robust aggregation rule, for $T \in \mathbb{N}^*$ iterations. We have for every $t \in \{1, \dots, T-1\}$:*

$$\mathbb{E}_{\mathcal{A}} \|\Sigma_{\mathcal{H},t}\|_{sp} = \mathbb{E}_{\mathcal{A}} [\lambda_{\max}(\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} (g_t^{(i)} - \bar{g}_t)(g_t^{(i)} - \bar{g}_t)^{\top})] \leq \sigma^2 + G^2.$$

Proof. In the following, we derive a bound for the stochastic case $\mathcal{A} = \text{SGD}$; the result naturally extends to the deterministic setting $\mathcal{A} = \text{GD}$ as well. The proof of the refined upper bound leverages the ability to compute the expectation of the controllable noise—specifically, the noise resulting from random sample selection $J_t^{(i)}$, $i \in \mathcal{H}$, $t \in \{0, \dots, T-1\}$ —while conditioning on the uncontrollable noise introduced by Byzantine workers. Since the bound is independent of the latter, we establish the result by applying the law of total expectation.

Let $t \in \{0, \dots, T-1\}$,

$$\Delta_t = \lambda_{\max}(\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} (g_t^{(i)} - \bar{g}_t)(g_t^{(i)} - \bar{g}_t)^{\top}) = \sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, g_t^{(i)} - \bar{g}_t \rangle^2.$$

We decompose,

$$g_t^{(i)} - \bar{g}_t = g_t^{(i)} - \nabla \widehat{R}_i(\theta_t) + \nabla \widehat{R}_i(\theta_t) - \nabla \widehat{R}_{\mathcal{H}}(\theta_t) + \nabla \widehat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t$$

so that we have:

$$\begin{aligned}
\langle v, g_t^{(i)} - \bar{g}_t \rangle^2 &= \langle v, \nabla \widehat{R}_i(\theta_t) - \nabla \widehat{R}_{\mathcal{H}}(\theta_t) \rangle^2 + \langle v, g_t^{(i)} - \nabla \widehat{R}_i(\theta_t) + \nabla \widehat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t \rangle^2 \\
&\quad + 2\langle v, g_t^{(i)} - \nabla \widehat{R}_i(\theta_t) + \nabla \widehat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t \rangle \langle v, \nabla \widehat{R}_i(\theta_t) - \nabla \widehat{R}_{\mathcal{H}}(\theta_t) \rangle.
\end{aligned}$$

Upon averaging, taking the supremum over the unit ball, and then taking total expectations, we get:

$$\begin{aligned}
\mathbb{E} \left[\sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, g_t^{(i)} - \bar{g}_t \rangle^2 \right] &\leq \mathbb{E} \left[\sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, \nabla \widehat{R}_i(\theta_t) - \nabla \widehat{R}_{\mathcal{H}}(\theta_t) \rangle^2 \right] \\
&\quad + \mathbb{E} \left[\sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, g_t^{(i)} - \nabla \widehat{R}_i(\theta_t) + \nabla \widehat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t \rangle^2 \right] \\
&\quad + 2\mathbb{E} \left[\sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, g_t^{(i)} - \nabla \widehat{R}_i(\theta_t) + \nabla \widehat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t \rangle \langle v, \nabla \widehat{R}_i(\theta_t) - \nabla \widehat{R}_{\mathcal{H}}(\theta_t) \rangle \right] \quad (21)
\end{aligned}$$

We show that the last term in eq. (21) is non-positive. In fact, with $M = N + N^{\top}$,

$$N = \sum_{i \in \mathcal{H}} (g_t^{(i)} - \nabla \widehat{R}_i(\theta_t) + \nabla \widehat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t)(\nabla \widehat{R}_i(\theta_t) - \nabla \widehat{R}_{\mathcal{H}}(\theta_t))^{\top}$$

and using Lemma D.4 from [5]:

$$\begin{aligned}
& \mathbb{E} \left[\sup_{\|v\|_2 \leq 1} 2 \sum_{i \in \mathcal{H}} \langle v, g_t^{(i)} - \nabla \hat{R}_i(\theta_t) + \nabla \hat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t \rangle \langle v, \nabla \hat{R}_i(\theta_t) - \nabla \hat{R}_{\mathcal{H}}(\theta_t) \rangle \right] \\
&= \mathbb{E} \left[\sup_{\|v\|_2 \leq 1} \langle v, Mv \rangle \right] \\
&\leq 9^d \sup_{\|v\|_2 \leq 1} \mathbb{E} [\langle v, Mv \rangle] \\
&= 9^d \sup_{\|v\|_2 \leq 1} 2 \mathbb{E} \sum_{i \in \mathcal{H}} \underbrace{\langle v, \mathbb{E}[g_t^{(i)} - \nabla \hat{R}_i(\theta_t) + \nabla \hat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t | \theta_t] \rangle}_{=0} \langle v, \nabla \hat{R}_i(\theta_t) - \nabla \hat{R}_{\mathcal{H}}(\theta_t) \rangle
\end{aligned}$$

We now bound the second term in eq. (21):

$$\begin{aligned}
& \mathbb{E} \left[\sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, g_t^{(i)} - \nabla \hat{R}_i(\theta_t) + \nabla \hat{R}_{\mathcal{H}}(\theta_t) - \bar{g}_t \rangle^2 \right] \\
&\leq \mathbb{E} \left[\sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, g_t^{(i)} - \nabla \hat{R}_i(\theta_t) \rangle^2 \right] \\
&\leq \mathbb{E} \left[\frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \underbrace{\mathbb{E}_{J_t^{(i)}} [\|g_t^{(i)} - \nabla \hat{R}_i(\theta_t)\|_2^2]}_{\leq \sigma^2} \right]
\end{aligned}$$

And finally we bound the first term in eq. (21) with Assumption F.1:

$$\begin{aligned}
& \mathbb{E} \sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, \nabla \hat{R}_i(\theta_t) - \nabla \hat{R}_{\mathcal{H}}(\theta_t) \rangle^2 \leq \sup_{\theta \in \Theta} \sup_{\|v\|_2 \leq 1} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \langle v, \nabla \hat{R}_i(\theta) - \nabla \hat{R}_{\mathcal{H}}(\theta) \rangle^2 \\
&= \sup_{\theta \in \Theta} \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \|\nabla \hat{R}_i(\theta) - \nabla \hat{R}_{\mathcal{H}}(\theta)\|_2^2 \leq G^2 \quad \square
\end{aligned}$$

G High-probability excess risk bounds for robust distributed learning under smooth and strongly convex loss functions

In this section, we establish high-probability excess risk bounds (see Equation (1)) under smooth and strongly convex assumptions for robust distributed learning algorithms, providing insights into their generalization performance. Specifically, we combine a deterministic optimization error bound—adapted directly from prior work [15]—with our high-probability stability bound.

High probability bound under uniform stability framework are derived as follows. The change of a single training sample has a small effect on a uniformly stable algorithm. This broadly implies low variance in the difference between empirical and population risk. If this difference also has low expectation, the empirical risk should accurately approximate the population risk [10]. In fact, we directly adapt the results from [31] to our framework and state the following high-probability bound.

Lemma G.1. *Consider the setting described in Section 2, under either Byzantine failure or poisoning attacks. Let \mathcal{A} an ε -uniformly stable. Assume $\forall z \in \mathcal{Z}, \ell(\cdot, z) \leq \ell_\infty$, then we have that for any $0 < \delta < 1$, with probability at least $1 - \delta$:*

$$\begin{aligned}
|R_{\mathcal{H}}(\mathcal{A}(\mathcal{S})) - \hat{R}_{\mathcal{H}}(\mathcal{A}(\mathcal{S}))| \in \mathcal{O} \left(\min \left\{ \varepsilon \ln(|\mathcal{H}|m) \ln\left(\frac{1}{\delta}\right) + \frac{\ell_\infty}{\sqrt{|\mathcal{H}|m}} \sqrt{\ln\left(\frac{1}{\delta}\right)}, \right. \right. \\
\left. \left. \left(\sqrt{\varepsilon \ell_\infty} + \frac{\ell_\infty}{\sqrt{|\mathcal{H}|m}} \right) \sqrt{\ln\left(\frac{1}{\delta}\right)} \right\} \right)
\end{aligned}$$

Poisoning attacks. By adapting the optimization bound technique from [15, Theorem 4] to GD algorithm with SMEA under poisoning attacks, we show that for a constant number of steps—specifically, when the step count satisfies $T \geq \ln(4\frac{f+1/m}{n-f} + \kappa) / \ln(1 - \frac{\mu}{L}) \in \mathcal{O}(1)$ —the expected excess risk is of the order of:

$$\frac{C^2}{\mu} \mathcal{O} \left(\frac{f}{n-f} + \frac{1}{m(n-f)} \right).$$

Additionally, we provide a high-probability bound on the excess risk using Lemma G.1. Specifically, for any $0 < \delta < 1$, with probability at least $1 - \delta$ the excess risk lies within:

$$\mathcal{O} \left(\min \left\{ \frac{f}{n-f} C^2 + \left(\frac{f}{n-f} + \frac{1}{m(n-f)} \right) C^2 \ln((n-f)m) \ln\left(\frac{1}{\delta}\right) + \frac{\ell_\infty}{\sqrt{(n-f)m}} \sqrt{\ln\left(\frac{1}{\delta}\right)}, \right. \right. \\ \left. \left. \frac{f}{n-f} C^2 + \left(C \sqrt{\ell_\infty \left(\frac{f}{n-f} + \frac{1}{m(n-f)} \right)} + \frac{\ell_\infty}{\sqrt{(n-f)m}} \right) \sqrt{\ln\left(\frac{1}{\delta}\right)} \right\} \right)$$

In the absence of a minimax statistical lower bound for the heterogeneous case—and acknowledging that the comparison is imperfect due to differing assumptions—it remains informative to compare our resulting bound with those in [47, 48], as well as with Proposition 1 in [4].

Byzantine failure attacks. Under Byzantine failure attacks, the derivation proceeds similarly and yields an analogous bound, with a dependence on $\sqrt{\kappa}$ replacing the $\frac{f}{n-f}$ term in the uniform stability bound.