

Robust Training in High Dimensions via Block Coordinate Geometric Median Descent

Anish Acharya, Abolfazl Hashemi, Prateek Jain, Sujay Sanghavi, Inderjit Dhillon, Ufuk Topcu



Robust DNN Training

- (Gross Corruption) Adversary can inspect all samples and replace $0 \leq \psi \leq 1/2$ fraction of them with **arbitrary** points. If \mathcal{G} and \mathcal{B} are sets of good and bad points $\alpha = \frac{|\mathcal{B}|}{|\mathcal{G}|} = \frac{\psi}{\psi-1} \leq 1$.
- The goal of this paper is to design an efficient first-order optimization method to solve *smooth non-convex optimization* problems with finite-sum structure (ERM formulation of DNN training), under gross corruption, *without any prior knowledge about the malicious samples*.

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) := \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} f_i(\mathbf{x})$$

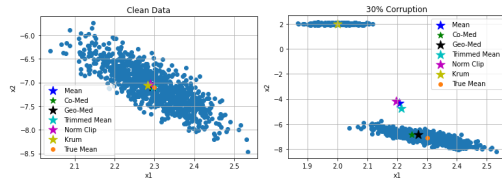
Vulnerability of SGD

- SGD: $\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \bar{\mathbf{g}}_t$, $\bar{\mathbf{g}}_t^{(i)} = \frac{1}{|\mathcal{D}_t|} \sum_{i \in \mathcal{D}_t} \nabla f_i(\mathbf{x}_t)$.
- Breakdown Point**: smallest fraction of contamination that must be introduced to cause an estimator to produce arbitrarily wrong estimates.
- A single corrupt sample can lead SGD to an *arbitrarily poor solution*. Consider a single malicious gradient: $\mathbf{g}_j^{(i)} = -\sum_{i \in \mathcal{D}_t \setminus j} \mathbf{g}_i^{(i)}$
- SGD has lowest possible *asymptotic breakdown of 0* under gross corruption due to the *linear gradient aggregation* step.

Robust SGD

- Replace Mean with *Robust Mean Estimator*.
- Geometric Median**: Optimal Breakdown point of $1/2$

$$\mathbf{x}_* = \text{GM}(\{\mathbf{x}_i\}) = \arg \min_{\mathbf{y} \in \mathbb{R}^d} \left[g(\mathbf{y}) := \sum_{i=1}^n \|\mathbf{y} - \mathbf{x}_i\| \right]$$

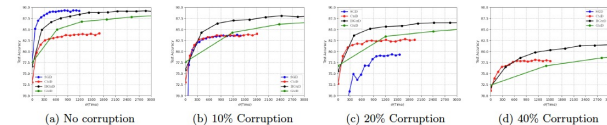


- GM Descent**: $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \hat{\mathbf{g}}_t$, $\hat{\mathbf{g}}_t = \text{GM}(\{\mathbf{g}_i\})$
- Finding ϵ -approximate GM of n points in \mathbb{R}^d requires at least $\mathcal{O}(d/\epsilon^2)$ compute making GM-SGD **computationally intractable** for DNN training e.g. $d \approx 60\text{M}$ Alexnet, $d \approx 175\text{B}$ GPT3

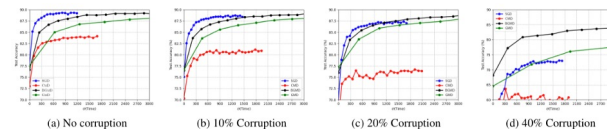
Block Geometric Median Descent

- gradient distribution of overparameterized DNNs are *long tailed* implying performing *gradient aggregation in low dimensions* should have *little impact* in the downstream optimization task.
- Judiciously *sample an informative block of k dimensions* to perform GM in \mathbb{R}^k ($k \ll d$)
- Keep track of Residual** information loss due to dimensionality reduction and add back to gradient estimate in future iterations.

Empirical Evidence



Feature Corruption Test accuracy as a *function of wall clock time* for training Fashion-MNIST using LeNet (1.16 M params) in presence of impulse noise.



Label Corruption. Test accuracy under label noise.

Theoretical Guarantees

Algorithm	Aggregation Operator*	Iteration Complexity [†]	Breakdown Point [‡]
SGD	MEAN(-)	$\mathcal{O}(bd)$	0
(Yang et al., 2019; Yin et al., 2018)	GM(-)	$\mathcal{O}(bd \log b)$	1/2
(Wu et al., 2020)	GM(-)	$\mathcal{O}(dc^{-2} + bd)$	1/2
BGMd (This work)	BGM(-)	$\mathcal{O}(kc^{-2} + bd)$	1/2
(Data and Diggavi, 2020)	(Steinhardt et al., 2017)	$\mathcal{O}(db^2 \min(d, b) + bd)$	1/4
(Blanchard et al., 2017)	KRUM(-)	$\mathcal{O}(b^2 d)$	$[\beta]$
(Yin et al., 2018)	CTM _{SL} (-)	$\mathcal{O}(bd(1-2\beta) + bd \log b)$	$[\beta]$
(Ghosh et al., 2019; Gupta et al., 2020)	N _{CB} (-)	$\mathcal{O}(bd(2-\beta) + b \log b)$	$[\beta]$

- Non-Convex and Smooth**: Suppose f_i corresponding to non-corrupt samples i.e. $i \in \mathcal{G}$ are *L smooth* and *non-convex*. Run BGMd with ϵ approx. GM oracle and $\gamma = \frac{1}{2L}$ in presence of α corruption for T iterations. Sample any iteration τ uniformly at random then:

$$\mathbb{E} \|\nabla f(\mathbf{x}_\tau)\|^2 = \mathcal{O} \left(\frac{LR_0}{T} + \frac{\sigma^2 \xi^{-2}}{(1-\alpha)^2} + \frac{L^2 \epsilon^2}{|\mathcal{G}|^2 (1-\alpha)^2} \right)$$

