

# Robust Training in High Dimensions via Block Coordinate Geometric Median Descent

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

## TL;DR

- While mini-batch SGD is the de-facto method to solve problems with **finite sum structure**; even a single grossly corrupted sample can lead SGD to an arbitrarily poor solution.
- The vulnerability of SGD is often attributed to the gradient aggregation using  $\text{MEAN}(\cdot)$  which has 0 breakdown point. This motivates replacing the gradient aggregation with robust estimators like geometric median (GM), which *achieves the optimal breakdown point* of  $\frac{1}{2}$ .
- Despite the strong robustness guarantees of GM, all the known methods to compute an approximate GM solution, scale poorly with the dimension of the problem ( $d$ ) making it prohibitive in high dimensional settings.
- We propose BGMD (Algorithm 1), a method for robust optimization in high dimensions, which is significantly more efficient than the standard GM-SGD method but is still able to maintain the optimal breakdown point.

**Definition 1 (Gross Corruption Model).** Given  $0 \leq \psi < \frac{1}{2}$  and a distribution family  $\mathcal{D}$  on  $\mathbb{R}^d$  the adversary operates as follows:  $n$  samples are drawn from  $D \in \mathcal{D}$ . The adversary is allowed to **inspect** all the samples and replace up to  $\psi n$  samples with arbitrary points. This implies that  $\alpha := |\mathbb{B}|/|\mathbb{G}| < 1$ , where  $\mathbb{B}$  and  $\mathbb{G}$  are the sets of corrupt and good samples.

**Definition 2 (Breakdown Point).** Breakdown point  $\eta$  of an estimator is the smallest fraction of contamination that must be introduced to cause an estimator to break e.g. an estimator has **optimal breakdown point**  $1/2$  if robust  $\forall \alpha < 1$  gross corruption.

**Definition 3 (Geometric Median).**

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

We call a point  $\mathbf{x} \in \mathbb{R}^d$  an  $\epsilon$ -accurate geometric median if  $g(\mathbf{x}) \leq (1 + \epsilon)g(\mathbf{x}_*)$ .

**(Overview of BGMD).** At a high level, at each iteration BGMD selects a block of  $0 < k \leq d$  important coordinates of the stochastic gradients using Algorithm 2. The remaining  $(d-k)$  dimensions are discarded and gradient aggregation happens only along these selected  $k$  directions. While descending along only a small subset of  $k$  coordinates at each iteration significantly improves the per iteration computational cost (Lemma 2), a smaller value of  $k$  would also imply larger gradient information loss a smaller  $\xi$  (Lemma 1). To mitigate this we propose a memory mechanism: Throughout training, we keep track of the residual error  $\|\mathbf{G}_t - \mathcal{C}_k(\mathbf{G}_t)\|$  incurred due to ignoring  $(d-k)$  dimensions via  $\hat{\mathbf{m}}_t \in \mathbb{R}^d$  that is appropriately added back in the subsequent iterations.

**Algorithm 1** Block GM Descent (BGMD)

**Initialize:** estimate:  $\mathbf{x}_0 \in \mathbb{R}^d$ , step-size:  $\gamma$ , memory:  $\hat{\mathbf{m}}_0 = \mathbf{0}$ , Block Coordinate Selection operator:  $\mathcal{C}_k(\cdot)$ , Geometric Median operator:  $\text{GM}(\cdot)$   
epochs  $t = 0, \dots$ , until convergence select samples  $\mathcal{D}_t = \{i_1, \dots, i_b\}$   
obtain:  $\mathbf{g}_t^{(i)} := \nabla f_i(\mathbf{x}_t)$ ,  $\forall i \in \mathcal{D}_t$  (back-propagation)  
Let  $\mathbf{G}_t \in \mathbb{R}^{b \times d}$  s.t. each row  $\mathbf{G}_t[i, :] = \mathbf{g}_t^{(i)}$   
 $\mathbf{G}_t[i, :] \leftarrow \gamma \mathbf{G}_t[i, :] + \hat{\mathbf{m}}_t \forall i \in [b]$  (add memory)  
 $\Delta_t := \mathcal{C}_k(\mathbf{G}_t) \in \mathbb{R}^{b \times k}$  (subset  $k$  dim via Algo. 2)  
 $\mathbf{M}_{t+1} = \mathbf{G}_t - \Delta_t$  (compute residuals)  
 $\hat{\mathbf{m}}_{t+1} = \frac{1}{b} \sum_{0 \leq i < b} \mathbf{M}_{t+1}[i, :]$  (update memory)  
 $\hat{\mathbf{g}}_t := \text{GM}(\Delta_t)$  (robust aggregation in  $\mathbb{R}^k$ )  
 $\mathbf{x}_{t+1} := \mathbf{x}_t - \hat{\mathbf{g}}_t$  (parameter update)

**Algorithm 2** Block Coordinate Selection Strategy

**Input:**  $\mathbf{G}_t \in \mathbb{R}^{n \times d}$ ,  $k$   
coordinates  $j = 0, \dots, d-1$   $s_j \leftarrow \|\mathbf{G}_t[:, j]\|^2$  (norm along each dimension) Sample set  $\mathbb{I}_k$  of  $k$  dimensions with probabilities proportional to  $s_j$   
 $\mathcal{C}_k(\mathbf{G}_t)[i, j \in \mathbb{I}_k] = \mathbf{G}_t[i, j]$ ,  $\mathcal{C}_k(\mathbf{G}_t)[i, j \notin \mathbb{I}_k] = 0$   
**Return:**  $\mathcal{C}_k(\mathbf{G}_t)$

## Convergence Guarantees

**Lemma 1 (Contraction Mapping).** Algorithm 2 yields a contraction approximation  $\mathbb{E} [\|\mathcal{C}_k(\mathbf{G}_t) - \mathbf{G}_t\|^2 | \mathbf{G}_t] \leq (1 - \xi)\|\mathbf{G}_t\|^2$ ,  $\frac{k}{d} \leq \xi \leq 1$ .

**Assumption 1 (Unbiased Oracle - Bounded Variance).**

$$\mathbb{E}_{z \sim \mathcal{D}_i} [\mathbf{g}_i(\mathbf{x}, z)] = \nabla f_i(\mathbf{x}) \quad (2)$$

$$\mathbb{E}_{z \sim \mathcal{D}_i} \|\nabla F_i(\mathbf{x}, z)\|^2 \leq \sigma^2 \quad (3)$$

**Assumption 2 (Smoothness).**

$$f_i(\mathbf{x}) \leq f_i(\mathbf{y}) + \langle \mathbf{x} - \mathbf{y}, \nabla f_i(\mathbf{y}) \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 \quad (4)$$

**Assumption 3 (Polyak-Lojasiewicz Condition).**

$$\|\nabla f(\mathbf{x})\|^2 \geq 2\mu(f(\mathbf{x}) - f(\mathbf{x}^*)), \mu > 0 \quad (5)$$

**Theorem 1 (Smooth Non-convex).** Suppose Assumption 1-2 hold. Run Algorithm 1 with compression factor  $\xi$  (Lemma 1), learning rate  $\gamma_t = 1/2L$  and  $\epsilon$ -approximate  $\text{GM}(\cdot)$  in presence of  $\alpha$ -corruption (Definition 1) for  $T$  iterations, then for any  $\tau \in [T]$  sampled uniformly at random:

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

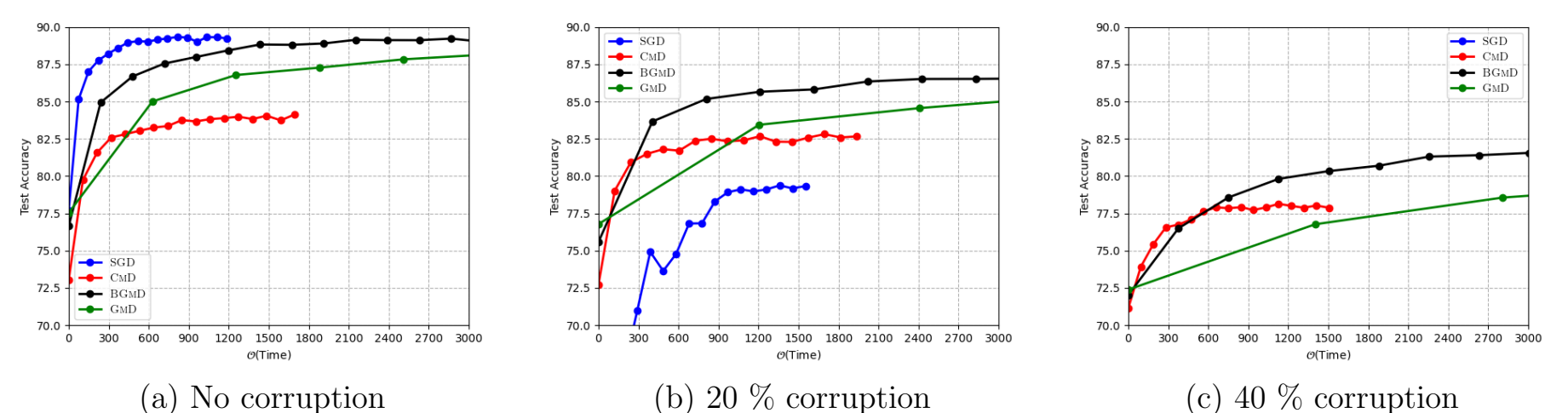
**Theorem 2 (Non-convex under PLC).** Suppose Assumption 1-3 hold. Then, after  $T$  iterations BGMD with compression factor  $\xi$ , learning rate  $\gamma_t = 1/4L$  and  $\epsilon$ -approximate  $\text{GM}(\cdot)$  oracle in presence of  $\alpha$ -corruption satisfies:

$$\mathbb{E} \|\hat{\mathbf{x}}_T - \mathbf{x}^*\|^2 = \mathcal{O} \left( \frac{LR_0}{\mu^2} \left[ 1 - \frac{\mu}{8L} \right]^T + \frac{\sigma^2 \xi^{-2}}{\mu^2 (1-\alpha)^2} + \frac{L^2 \epsilon^2}{\mu^2 |\mathbb{G}|^2 (1-\alpha)^2} \right)$$

where,  $\hat{\mathbf{x}}_T := \frac{1}{W_T} \sum_{t=0}^{T-1} w_t \mathbf{x}_t$ ,  $W_T := \sum_{t=0}^{T-1} w_t$  with weights  $w_t := (1 - \frac{\mu}{8L})^{-(t+1)}$ .

## Empirical Evidence

**Lemma 1 (Linear Speedup).** For,  $\beta \leq \mathcal{O}(1/F - b\epsilon^2)$ , given an  $\epsilon$ -approximate GM oracle, Algorithm 1 achieves a factor  $F$  speedup over GM-SGD.



Robustness to Feature Corruption (Accuracy over clock time)

We observe with extensive experiments under different sources and models of corruption - Feature Corruption; Gradient Corruption; Label Corruption; GM based methods are indeed superior while standard SGD or CMD can be significantly inaccurate. By judiciously choosing  $k$ , BGMD can be more efficient than GM, often resulting in more than 3x speedup. Despite using small  $\beta \leq 0.15$  in all our experiments, it retains high generalization performance emphasizing the role of memory mechanism.