

Centered Clipping with Per-client Thresholds given Side Information

1 Byzantine Impact Bound via Side Information

We give a bound on the *aggregate Byzantine impact* that holds under omniscient attackers. The bound depends only on quantities that are either (i) directly controlled by side information and optimisation tolerance, or (ii) intrinsic to the honest cohort.

Setup and Notation. Let \mathcal{G} and \mathcal{B} be the honest and Byzantine (client) index sets, with $|\mathcal{B}| = \delta n$ and $|\mathcal{G}| = (1 - \delta)n$. At the current round, the initial guess for the aggregation is $\mathbf{x}_0 \in \mathbb{R}^d$ and client proposals are $\mathbf{x}_i \in \mathbb{R}^d$. Define

$$\mathbf{d}_i := \frac{\mathbf{x}_i - \mathbf{x}_0}{\|\mathbf{x}_i - \mathbf{x}_0\|} \quad (\mathbf{x}_i \neq \mathbf{x}_0), \quad \alpha_i(\boldsymbol{\nu}) := \min\left(1, \frac{\nu_i}{\|\mathbf{x}_i - \mathbf{x}_0\|}\right) \in [0, 1].$$

The one-step clipped aggregate is

$$\hat{\mathbf{x}}_{\boldsymbol{\nu}} = \mathbf{x}_0 + \frac{1}{n} \sum_{i=1}^n \alpha_i(\boldsymbol{\nu}) (\mathbf{x}_i - \mathbf{x}_0).$$

Let $\mathbf{g}_{\mathcal{V}}$ be the validation gradient (side information). Assume the $\boldsymbol{\nu}$ -selection is solved up to tolerance ε_{ν} such that

$$\|\mathbf{g}_{\mathcal{V}} - \hat{\mathbf{x}}_{\boldsymbol{\nu}}\| \leq \varepsilon_{\nu}. \quad (1)$$

Let the honest mean be $\bar{\mathbf{x}} := 1/|\mathcal{G}| \sum_{i \in \mathcal{G}} \mathbf{x}_i$ and define the validation bias w.r.t. the honest mean

$$\varepsilon_{\mathcal{V}} := \|\mathbf{g}_{\mathcal{V}} - \bar{\mathbf{x}}\|.$$

We also define the (average) honest dispersion

$$\bar{\zeta}_h := \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} \|\mathbf{x}_i - \bar{\mathbf{x}}\|.$$

Byzantine Aggregate. We denote the aggregate Byzantine contribution by

$$\mathbf{B} := \frac{1}{n} \sum_{j \in \mathcal{B}} \alpha_j(\boldsymbol{\nu}) (\mathbf{x}_j - \mathbf{x}_0), \quad \text{so that} \quad \hat{\mathbf{x}}_{\boldsymbol{\nu}} - \mathbf{x}_0 = \mathbf{B} + \underbrace{\frac{1}{n} \sum_{i \in \mathcal{G}} \alpha_i(\boldsymbol{\nu}) (\mathbf{x}_i - \mathbf{x}_0)}_{=: \mathbf{G}}.$$

Theorem 1.1 (Byzantine impact bound given side information). *Under the setup above, for arbitrary Byzantine choices $\{\mathbf{x}_j\}_{j \in \mathcal{B}} \subset \mathbb{R}^d$ and the one-step $\boldsymbol{\nu}$ selected to satisfy (1), the Byzantine aggregate obeys*

$$\|\mathbf{B}\| \leq (2 - \delta) \|\mathbf{g}_{\mathcal{V}} - \mathbf{x}_0\| + \varepsilon_{\nu} + (1 - \delta) (\varepsilon_{\mathcal{V}} + \bar{\zeta}_h). \quad (2)$$

Consequently, for any $\delta \in [0, 1]$, the Byzantine impact can be made arbitrarily small by driving the controllable quantities $\|\mathbf{g}_{\mathcal{V}} - \mathbf{x}_0\|$, ε_{ν} , and $\varepsilon_{\mathcal{V}}$ via curation of validation dataset and tighter optimisation for the robust aggregation.

Remark 1.2 (Behaviour as $\delta \rightarrow 1$). When $\delta \rightarrow 1$ (i.e., $|\mathcal{G}| \rightarrow 0$), the honest term vanishes and the bound reduces to

$$\|\mathbf{B}\| \leq \|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\| + \varepsilon_\nu.$$

Formally, $\bar{\mathbf{x}}$ and $\bar{\zeta}_h$ are defined only when $|\mathcal{G}| \geq 1$; the display should be read as the $\delta \rightarrow 1$ *limit* of (2). In words: even if the attacker set occupies (almost) all clients, the Byzantine impact remains controlled by *side-information alignment* $\|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\|$ and the *optimisation tolerance* ε_ν .

Proof. Starting from the decomposition $\mathbf{g}_\mathcal{V} - \mathbf{x}_0 = (\mathbf{g}_\mathcal{V} - \hat{\mathbf{x}}_\nu) + (\hat{\mathbf{x}}_\nu - \mathbf{x}_0) = \mathbf{e} + (\mathbf{B} + \mathbf{G})$, where $\mathbf{e} := \mathbf{g}_\mathcal{V} - \hat{\mathbf{x}}_\nu$, we have

$$\mathbf{B} = \mathbf{g}_\mathcal{V} - \mathbf{x}_0 - \mathbf{e} - \mathbf{G}.$$

Taking norms and using (1),

$$\|\mathbf{B}\| \leq \|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\| + \|\mathbf{e}\| + \|\mathbf{G}\| \leq \|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\| + \varepsilon_\nu + \|\mathbf{G}\|. \quad (3)$$

Since $\alpha_i(\nu) \in [0, 1]$, we have

$$\|\mathbf{G}\| = \left\| \frac{1}{n} \sum_{i \in \mathcal{G}} \alpha_i(\nu) (\mathbf{x}_i - \mathbf{x}_0) \right\| \leq \frac{1}{n} \sum_{i \in \mathcal{G}} \alpha_i(\nu) \|\mathbf{x}_i - \mathbf{x}_0\| \leq \frac{1}{n} \sum_{i \in \mathcal{G}} \|\mathbf{x}_i - \mathbf{x}_0\|,$$

which substituted into (3) yields

$$\|\mathbf{B}\| \leq \|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\| + \varepsilon_\nu + \frac{1}{n} \sum_{i \in \mathcal{G}} \|\mathbf{x}_i - \mathbf{x}_0\|.$$

For (2), observe that

$$\frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} \|\mathbf{x}_i - \mathbf{x}_0\| \leq \frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} (\|\mathbf{x}_i - \bar{\mathbf{x}}\| + \|\bar{\mathbf{x}} - \mathbf{x}_0\|) = \bar{\zeta}_h + \|\bar{\mathbf{x}} - \mathbf{x}_0\|.$$

Moreover,

$$\|\bar{\mathbf{x}} - \mathbf{x}_0\| \leq \|\bar{\mathbf{x}} - \mathbf{g}_\mathcal{V}\| + \|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\| = \varepsilon_\nu + \|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\|.$$

Combining the two displays gives

$$\frac{1}{|\mathcal{G}|} \sum_{i \in \mathcal{G}} \|\mathbf{x}_i - \mathbf{x}_0\| \leq \bar{\zeta}_h + \varepsilon_\nu + \|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\|.$$

Multiplying by $|\mathcal{G}|/n = (1 - \delta)$ and substituting into (3) yields

$$\|\mathbf{B}\| \leq \|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\| + \varepsilon_\nu + (1 - \delta) (\bar{\zeta}_h + \varepsilon_\nu + \|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\|) = (2 - \delta) \|\mathbf{g}_\mathcal{V} - \mathbf{x}_0\| + \varepsilon_\nu + (1 - \delta)(\varepsilon_\nu + \bar{\zeta}_h),$$

which is (2). \square