

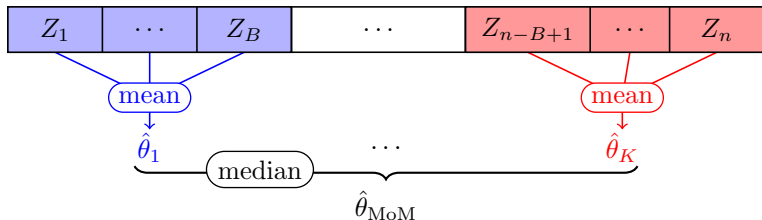
Generalization Bounds in the Presence of Outliers: a Median-of-Means Study

Pierre Laforgue¹, Guillaume Staerman², Stephan Cléménçon²

¹ Università degli Studi di Milano, Italy

² LTCI, Télécom Paris, Institut Polytechnique de Paris, France

The Median-of-Means (MoM) for heavy-tailed data



Z_1, \dots, Z_n i.i.d. realizations of r.v. Z s.t. $\mathbb{E}[Z] = \theta$, $\text{Var}(Z) = \sigma^2$.

$\forall \delta \in [e^{1-\frac{2n}{9}}, 1[$, for $K = \lceil \frac{9}{2} \ln(1/\delta) \rceil$ it holds [Devroye et al. 2016]:

$$\mathbb{P} \left\{ |\hat{\theta}_{\text{MoM}} - \theta| > 3\sqrt{6}\sigma \sqrt{\frac{1 + \ln(1/\delta)}{n}} \right\} \leq \delta.$$

The Median-of-Means (MoM) for outliers

$\{Z_1, \dots, Z_n\}$ contains $n - n_0$ *inliers* drawn i.i.d. from P , and n_0 *outliers*. We denote $\varepsilon = n_0/n$. Choosing $K = \lceil \beta(\varepsilon) \log(1/\delta) \rceil$, we have w.p.a.l. $1 - \delta$:

$$|\hat{\theta}_{\text{MoM}} - \theta| \leq \frac{12\sqrt{5}e\sigma}{(1-2\varepsilon)^{3/2}} \sqrt{\frac{1 + \log(1/\delta)}{n}}.$$

If in addition P is ρ sub-Gaussian, with $K = \lceil \alpha(\varepsilon)n \rceil$, we have w.p.a.l. $1 - \delta$:

$$|\hat{\theta}_{\text{MoM}} - \theta| \leq \frac{4\sqrt{5}\rho}{\sqrt{1-2\varepsilon}} \sqrt{\frac{\log(1/\delta)}{n}}.$$

If furthermore $n_0 \leq C_0 n^{\alpha_0}$, with the same K we have:

$$\mathbb{E} [|\hat{\theta}_{\text{MoM}} - \theta|] \leq \frac{2\sqrt{5}\rho}{\sqrt{1-2\varepsilon}} \left(4C_0 \frac{\Delta(\varepsilon)}{n^{(1-\alpha_0)/2}} + \sqrt{\frac{\pi}{n}} \right).$$

Similar guarantees for U -statistics, with application to Integral Probability Metrics [Staerman et al. 2021]

Generalization bounds for pairwise learning

MoU minimization (adaptation fom [Lecué et al. 2018]):

$$\hat{g}_{MoU} = \operatorname{argmin}_{g \in \mathcal{G}} \operatorname{median} \left(\sum_{i < j \in \mathcal{B}_1} \ell(g, Z_i, Z_j), \dots, \sum_{i < j \in \mathcal{B}_K} \ell(g, Z_i, Z_j) \right).$$

Algorithm 1 MoU Gradient Descent (MoU-GD)

input: $S_n, K, T \in \mathbb{N}^*, (\gamma_t)_{t \leq T} \in \mathbb{R}_+^T, u_0 \in \mathbb{R}^p$

for epoch from 1 to T **do**

```
// Randomly partition the data
Choose a random permutation  $\pi$  of  $\{1, \dots, n\}$ 
Build a partition  $B_1, \dots, B_k$  of  $\{\pi(1), \dots, \pi(n)\}$ 
// Select block with median risk
for  $k \leq K$  do
|  $\hat{U}_{B_k} = \sum_{i < j \in B_k^2} \ell(g_{u_t}, Z_i, Z_j)$ 
Set  $B_{\text{med}}$  s.t.  $\hat{U}_{B_{\text{med}}} = \operatorname{median}(\hat{U}_{B_1}, \dots, \hat{U}_{B_K})$ 
// Gradient step
 $u_{t+1} = u_t - \gamma_t \sum_{i < j \in B_{\text{med}}^2} \nabla_{u_t} \ell(g_{u_t}, Z_i, Z_j)$ 
```

return u_T

With guarantees in the presence of outliers:

$$\mathcal{R}(\hat{g}_{\text{alg}}) - \mathcal{R}(g^*) \leq \frac{8\sqrt{10M}}{\sqrt{1-2\epsilon}} \sqrt{\frac{\text{VC}_{\dim}(\mathcal{G})(1 + \log(n)) + \log(1/\delta)}{n}}.$$

Numerical experiments

Application to metric learning on the *iris* dataset:

