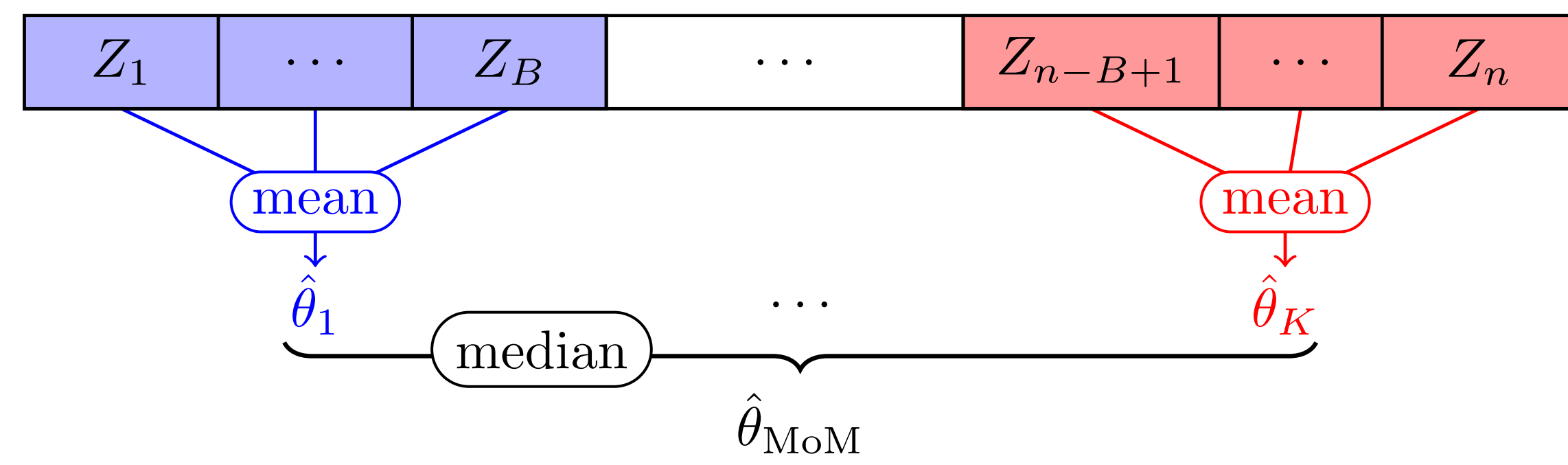


# Generalization Bounds in the Presence of Outliers: a Median-of-Means Study

P. Laforgue<sup>1</sup>, G. Staerman<sup>2</sup>, S. Cléménçon<sup>2</sup>

<sup>1</sup> University of Milan, Italy <sup>2</sup> LTCI, Télécom ParisTech, Institut Polytechnique de Paris, France

## MoM in the presence of outliers



**Assumption 1.** The sample  $\mathcal{S}_n = \{Z_1, \dots, Z_n\}$  contains  $n - n_O$  inliers drawn i.i.d. from  $P$ , and  $n_O$  outliers, upon which no assumption is made. Let  $\varepsilon = n_O/n$  be the fraction of outliers among sample  $\mathcal{S}_n$ .

**Assumption 2.** Let  $\alpha: [0, 1/2] \rightarrow [0, 1]$  be such that:  $\forall \varepsilon \in (0, 1/2)$ ,  $2\varepsilon < \alpha(\varepsilon) < 1$ . From mapping  $\alpha$ , we define the following functions:

$$\beta: \varepsilon \mapsto \frac{2\alpha(\varepsilon)}{\alpha(\varepsilon) - 2\varepsilon}, \quad \gamma: \varepsilon \mapsto \frac{\sqrt{\alpha(\varepsilon)(\alpha(\varepsilon) - \varepsilon)}}{(\alpha(\varepsilon) - 2\varepsilon)^{3/2}},$$

$$\Gamma: \varepsilon \mapsto \sqrt{\frac{\alpha(\varepsilon)}{\alpha(\varepsilon) - 2\varepsilon}}, \quad \Delta: \varepsilon \mapsto \sqrt{\frac{\alpha(\varepsilon)}{\varepsilon}}.$$

**Proposition 1.** Let  $\mathcal{S}_n$  and  $\alpha, \beta, \gamma, \Gamma, \Delta$  satisfying Assumptions 1 and 2 respectively. Then, for any  $\delta \in [e^{-n/\beta(\varepsilon)}, e^{-n\alpha(\varepsilon)/\beta(\varepsilon)}]$ , choosing  $K = \lceil \beta(\varepsilon) \log(1/\delta) \rceil$ , it holds w.p.a.l.  $1 - \delta$ :

$$|\hat{\theta}_{\text{MoM}} - \theta| \leq 4\sqrt{e}\sigma \gamma(\varepsilon) \sqrt{\frac{1 + \log(1/\delta)}{n}}.$$

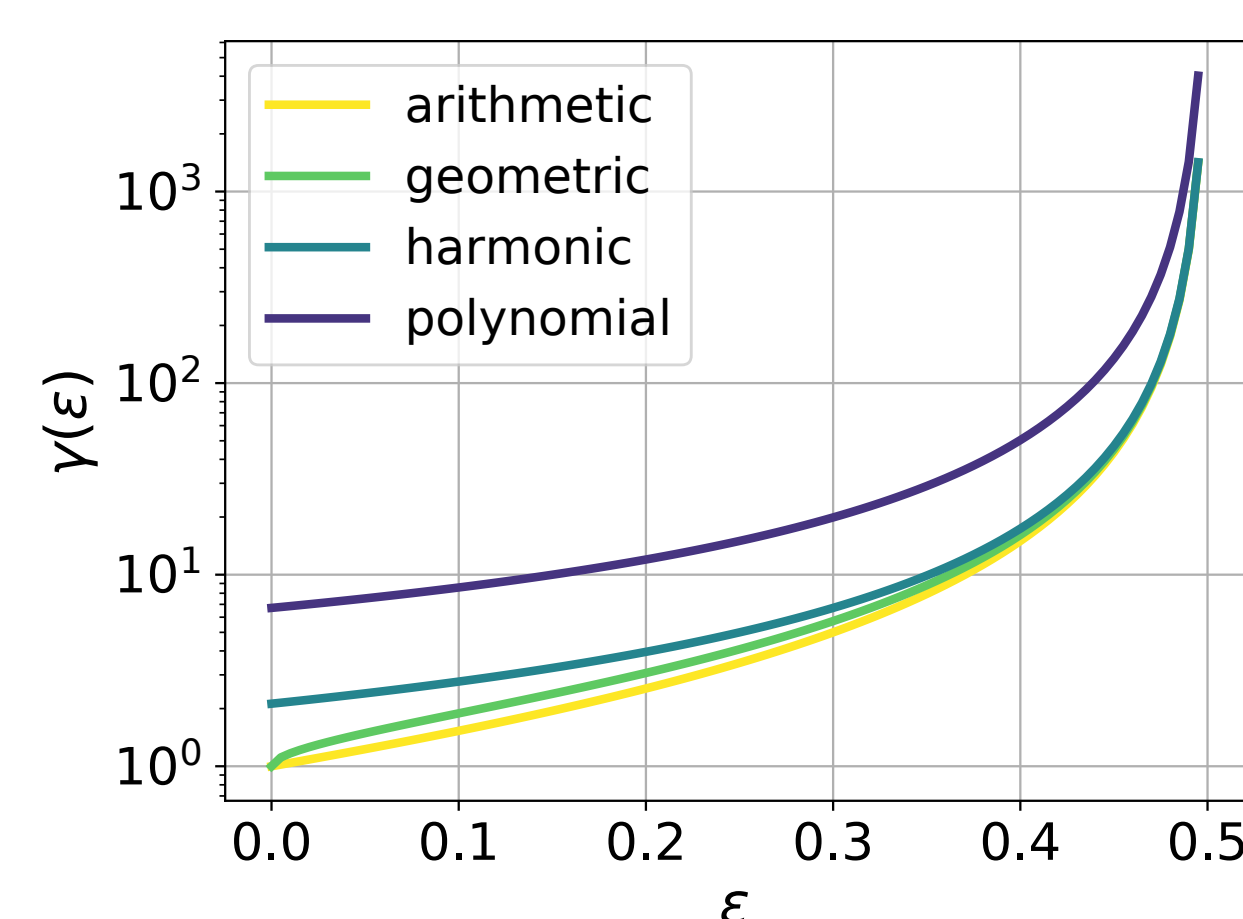
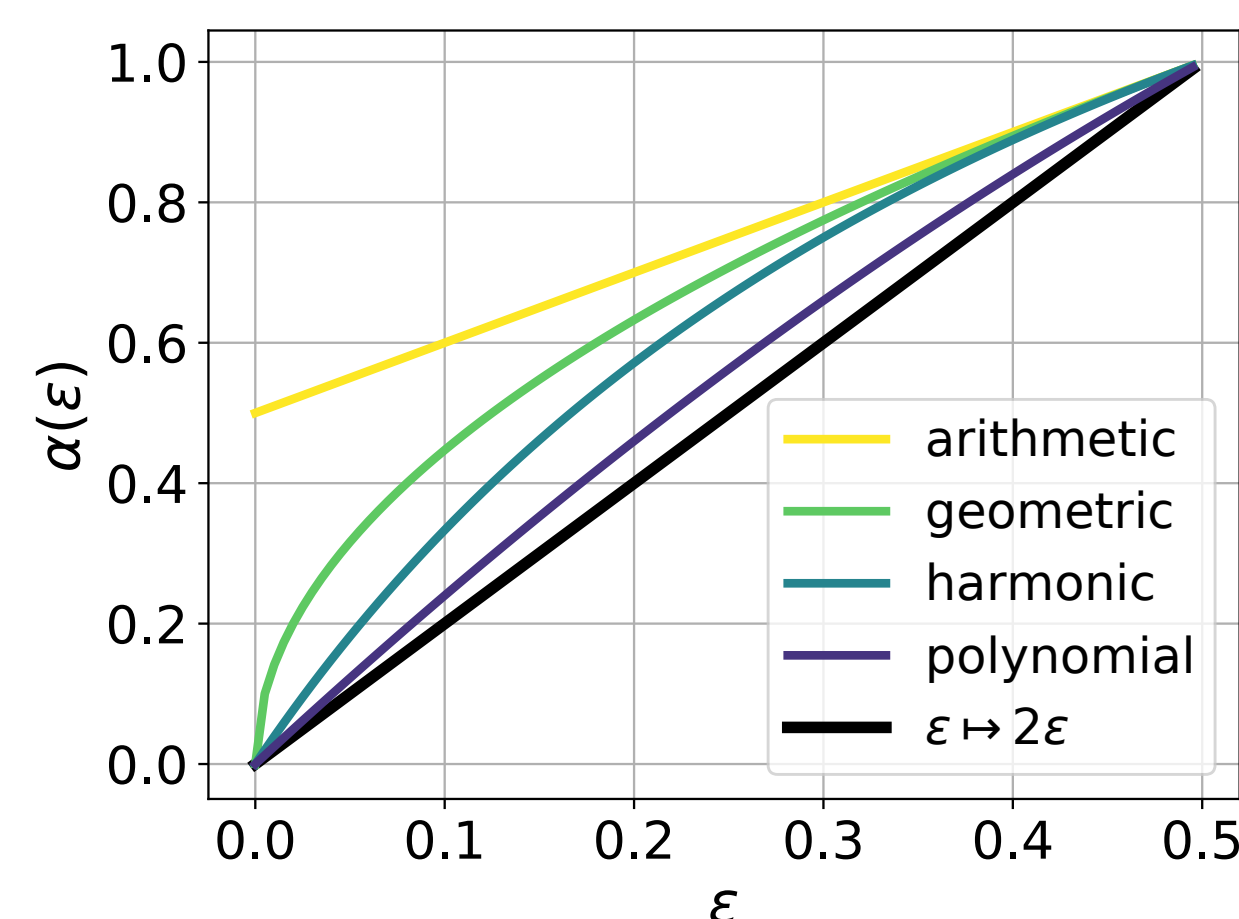
If in addition distribution  $P$  is  $\rho$  sub-Gaussian, then for all  $\delta \in (0, e^{-4n\alpha(\varepsilon)})$ , with  $K = \lceil \alpha(\varepsilon)n \rceil$ , it holds w.p.a.l.  $1 - \delta$ :

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

If furthermore  $n_O \leq C_0 n^{\alpha_0}$ , the same  $K$  gives:

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

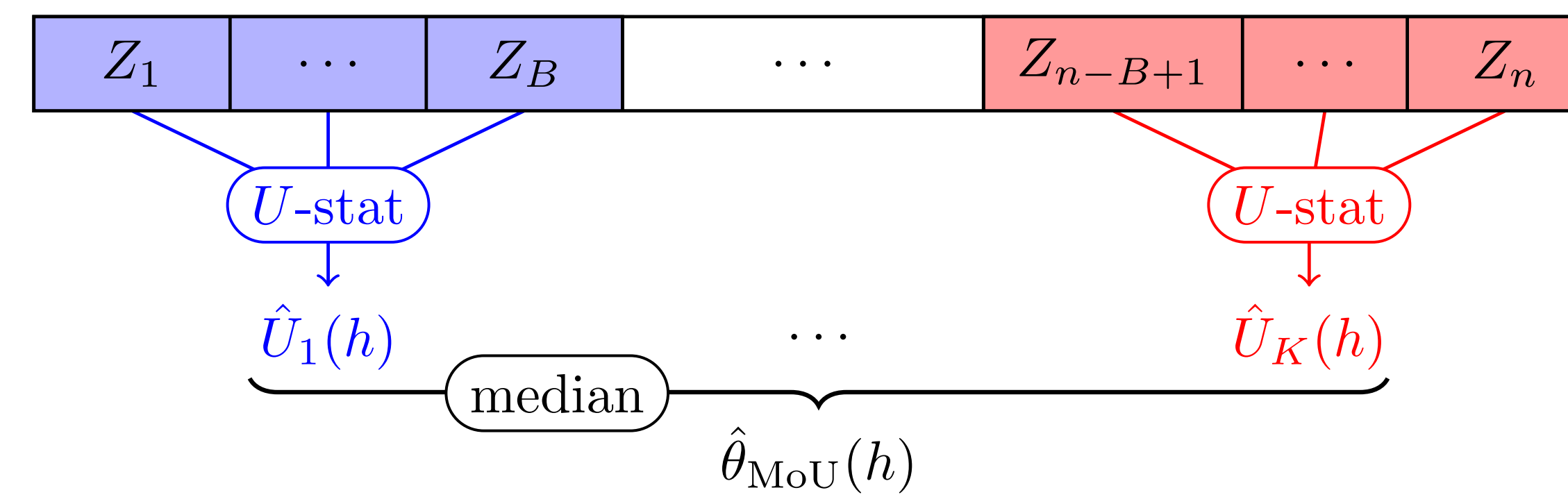
For most  $\alpha$ , we have  $\gamma(\varepsilon) \leq 3\sqrt{5}/(1-2\varepsilon)^{3/2}$  and  $\Gamma(\varepsilon) \leq \sqrt{5}/\sqrt{1-2\varepsilon}$ .



## MoU in the presence of outliers

A  $U$ -statistic is the MVU estimator of  $\mathbb{E}[h(Z_1, \dots, Z_d)]$ . It is given by

$$\hat{U}_n(h) = \frac{1}{\binom{n}{d}} \sum_{1 \leq i_1 < \dots < i_d \leq n} h(Z_{i_1}, \dots, Z_{i_d}).$$

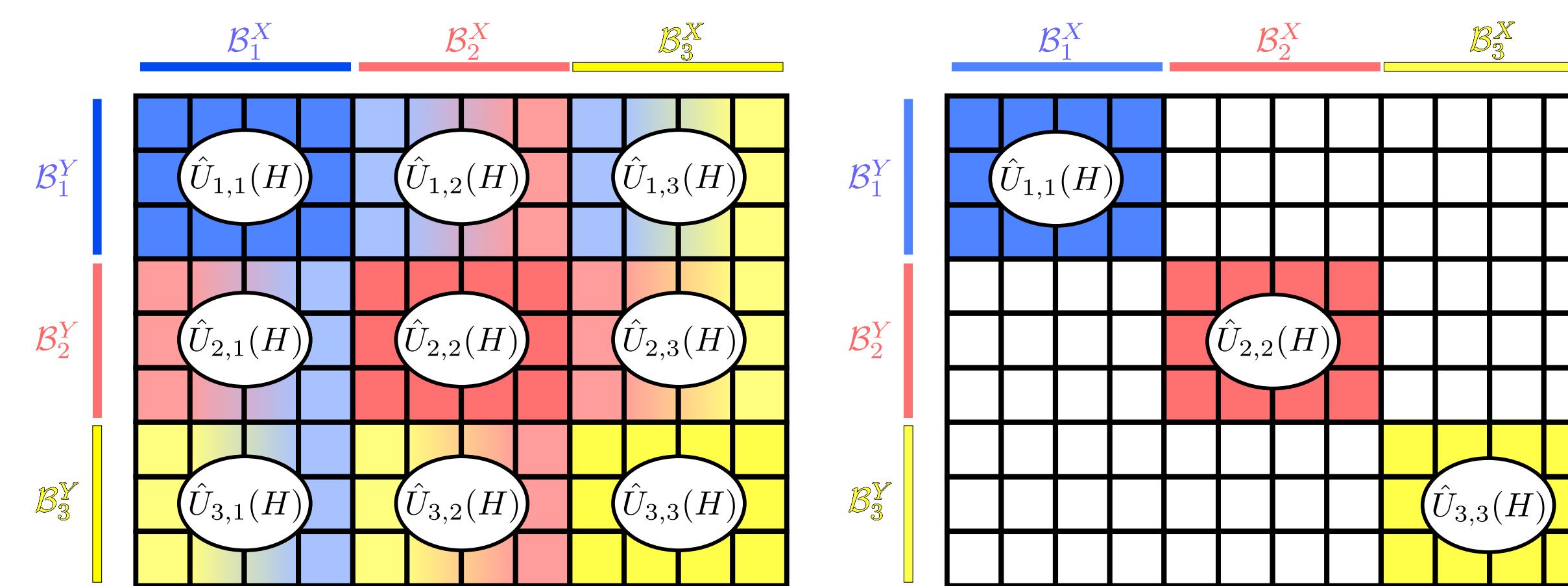


**Proposition 2.** Let  $\Sigma^2(h) = d! \sum_{c=1}^d \binom{d}{c} \zeta_c(h)$ , and  $M = \|h(Z_1, \dots, Z_d)\|_\infty$ . Then, under the assumptions of Proposition 1, with the same choices of  $K$  (and the same ranges of  $\delta$ ), we have w.p.a.l.  $1 - \delta$ :

$$|\hat{\theta}_{\text{MoU}}(h) - \theta(h)| \leq 4\sqrt{e} \Sigma(h) \gamma(\varepsilon) \sqrt{\frac{1 + \log(1/\delta)}{n}},$$

$$|\hat{\theta}_{\text{MoU}}(h) - \theta(h)| \leq 4\sqrt{d} M \Gamma(\varepsilon) \sqrt{\frac{\log(1/\delta)}{n}},$$

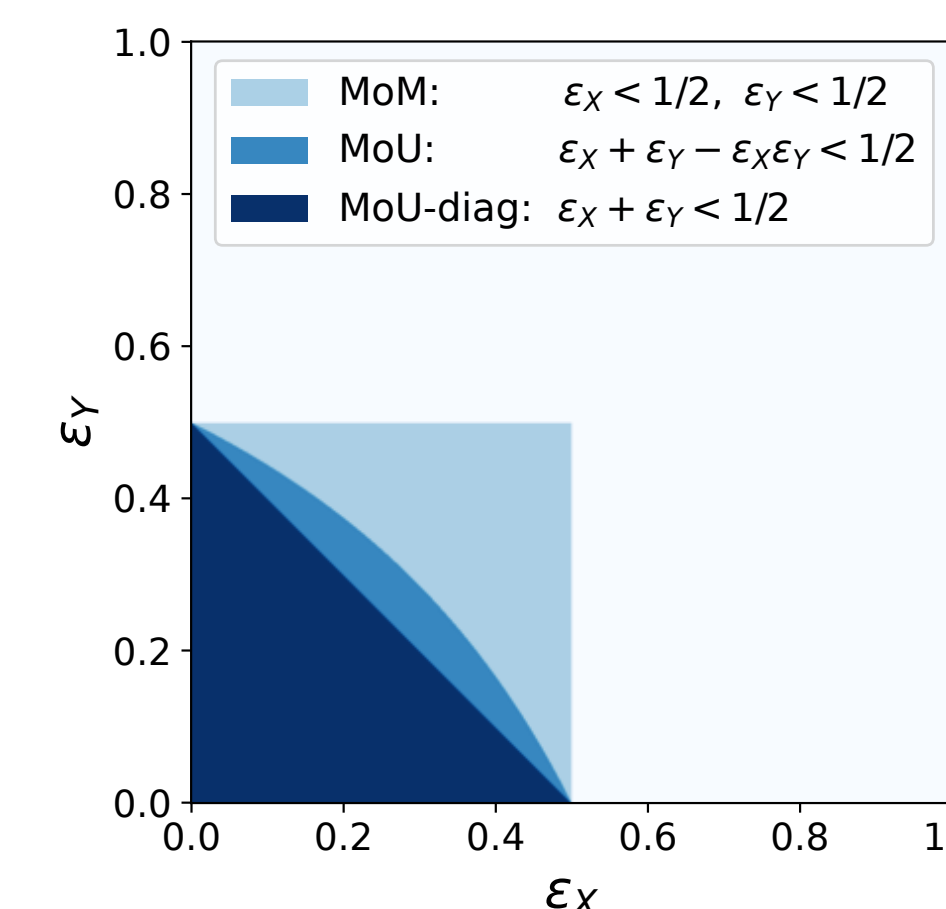
$$\mathbb{E} \left[ |\hat{\theta}_{\text{MoU}}(h) - \theta(h)| \right] \leq 2\sqrt{d} M \Gamma(\varepsilon) \left( 4C_0 \frac{\Delta(\varepsilon)}{n^{(1-\alpha_0)/2}} + \sqrt{\frac{\pi}{n}} \right).$$



The 2-sample  $U$ -statistic of degrees  $(1, 1)$  is defined as  $\hat{U}_{n,m}(H) = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m H(X_i, Y_j)$ . For IPMs,  $h(X, Y) = \sup_\phi \phi(X) - \phi(Y)$ .

Results are similar to the ones in Propositions 1 and 2, but with the following constants, and the notation  $\tilde{\varepsilon} = \varepsilon_X + \varepsilon_Y - \varepsilon_X \varepsilon_Y$ .

Asm.	Complete	Diag.
None	$\gamma(\tilde{\varepsilon})$	$\gamma(\varepsilon_X + \varepsilon_Y)$
$\ H\ _\infty < +\infty$	$\emptyset$	$\Gamma(\varepsilon_X + \varepsilon_Y)$
$n_O = \mathcal{O}(n^{\alpha_0})$	$\emptyset$	$\Delta(\varepsilon_X + \varepsilon_Y)$



## Outlier robust pairwise learning

**Goal:**  $g^* = \operatorname{argmin}_{g \in \mathcal{G}} \left\{ \mathcal{R}(g) = \mathbb{E}[\ell_g(Z, Z')] \right\}$  (ranking, metric learning)

**Approach:**  $\hat{g}_{\text{MoU}} = \operatorname{argmin}_{g \in \mathcal{G}} \operatorname{median} \left( \sum_{i < j \in \mathcal{B}_k} \ell_g(Z_i, Z_j), \text{ for } k \leq K \right)$

**Algorithm:** adapted from Lecué et al. 2018

### Algorithm 1 MoU Gradient Descent (MoU-GD)

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

**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**

$\hat{U}_{B_k} = \sum_{i < j \in B_k} \ell(g_{u_t}, Z_i, Z_j)$

    Set  $B_{\text{med}}$  s.t.  $\hat{U}_{B_{\text{med}}} = \operatorname{median}(\hat{U}_{B_1}, \dots, \hat{U}_{B_K})$

    // Gradient step

$u_{t+1} = u_t - \gamma_t \sum_{i < j \in B_{\text{med}}} \nabla_{u_t} \ell(g_{u_t}, Z_i, Z_j)$

**return**  $u_T$

**Theorem 1.** Suppose that  $\ell_g(Z, Z') < M$ , and that  $\mathcal{G}$  has finite VC-dimension. Under some technical assumptions taken from Lecué et al. 2018, the output of Algorithm 1 run on a corrupted sample  $\mathcal{S}_n$  converges almost surely towards  $\hat{g}_{\text{alg}}$ , that satisfies w.p.a.l.  $1 - \delta$ :

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

