

# Confidence regions and minimax rates in outlier-robust estimation on the probability simplex

Amir-Hossein Bateni

amirhossein.bateni@ensae.fr  
ENSAE, IP Paris

## Introduction

Assume  $X_1, \dots, X_n$  are  $n$  i.i.d. random variables taking their values in the  $k$ -dimensional probability simplex

$$\Delta^{k-1} = \{v \in \mathbb{R}_+^k : v_1 + \dots + v_k = 1\}.$$

Our goal is to estimate the unknown vector

$$\theta = \mathbf{E}[X_i]$$

in the case where the observations are contaminated by outliers.

Particular case: The distribution of a discrete random variable  $X$  taking  $k$  distinct values ( $X_i$ 's take values in  $\{e_1, \dots, e_k\}$ ).

## Minimax risk

We study the estimation error under three different metrics: total-variation, Hellinger and  $\mathbb{L}^2$  distances

$$\begin{aligned} d_{TV}(\hat{\theta}, \theta) &:= 1/2 \|\hat{\theta} - \theta\|_1, \\ d_H(\hat{\theta}, \theta) &:= 1/\sqrt{2} \|\hat{\theta}^{1/2} - \theta^{1/2}\|_2, \\ d_{\mathbb{L}^2}(\hat{\theta}, \theta) &:= \|\hat{\theta} - \theta\|_2, \end{aligned}$$

and we evaluate the minimax risk

$$\mathfrak{R}_{\square}(n, k, s, \varepsilon) := \inf_{\bar{\theta}_n} \sup_{P, Q} \mathbf{E}[d_{\square}(\bar{\theta}_n, \theta)],$$

where the *inf* is over all estimators  $\bar{\theta}_n$  built upon the contaminated observations and the *sup* is over all distributions  $P, Q$  on the probability simplex such that the mean  $\theta$  of  $P$  is  $s$ -sparse. The subscript  $\square$  of  $\mathfrak{R}$  above refers to the distance used in the risk, so that  $\square$  is TV, H, or  $\mathbb{L}^2$ .

## Lower and upper bounds

To ease notation, we let  $R_{\square}^{\square}(n, \varepsilon, \Theta, \hat{\theta})$  to be the worst-case risk of an estimator  $\hat{\theta}$ , where  $\square$  is either HC or AC. More precisely, for  $\mathcal{M}_n^{\square}(\varepsilon, \Theta) := \cup_{\theta \in \Theta} \mathcal{M}_n^{\square}(\varepsilon, \theta)$ , we set

$$R_{\square}^{\square}(n, \varepsilon, \Theta, \hat{\theta}) := \sup_{P_n \in \mathcal{M}_n^{\square}(\varepsilon, \Theta)} \mathbf{E}[d(\hat{\theta}, \theta^*)].$$

We denote by  $\Delta_s^{k-1}$  the set of all  $v \in \Delta^{k-1}$  having at most  $s$  non-zero entries.

**Theorem 1.** For every triple of positive integers  $(k, s, n)$  and for every  $\varepsilon \in [0, 1]$ , the sample mean  $\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i$  satisfies

$$\begin{aligned} R_{TV}^{AC}(n, \varepsilon, \Delta_s^{k-1}, \bar{X}_n) &\leq (s/n)^{1/2} + 2\varepsilon, \\ R_H^{AC}(n, \varepsilon, \Delta_s^{k-1}, \bar{X}_n) &\leq (s/n)^{1/2} + \sqrt{2}\varepsilon^{1/2}, \\ R_{\mathbb{L}^2}^{AC}(n, \varepsilon, \Delta_s^{k-1}, \bar{X}_n) &\leq (1/n)^{1/2} + \sqrt{2}\varepsilon. \end{aligned}$$

**Theorem 2.** There are universal constants  $c > 0$  and  $n_0$ , such that for any integers  $k \geq 3, s \leq k \wedge n, n \geq n_0$  and for any  $\varepsilon \in [0, 1]$ , we have

$$\begin{aligned} \inf_{\bar{\theta}_n} R_{TV}^{HC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq c\{(s/n)^{1/2} + \varepsilon\}, \\ \inf_{\bar{\theta}_n} R_H^{HC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq c\{(s/n)^{1/2} + \varepsilon^{1/2}\}, \\ \inf_{\bar{\theta}_n} R_{\mathbb{L}^2}^{HC}(n, \varepsilon, \Delta_s^{k-1}, \bar{\theta}_n) &\geq c\{(1/n)^{1/2} + \varepsilon\}, \end{aligned}$$

where  $\inf_{\bar{\theta}_n}$  stands for the infimum over all measurable functions  $\bar{\theta}_n$  from  $(\Delta^{k-1})^n$  to  $\Delta^{k-1}$ .

Thus, the rate obtained for the sample mean is minimax optimal.

## Instance based bounds

**Theorem 4.** Suppose  $X_i$ 's take value in  $\{e_1, e_2, \dots\}$ , and for  $j \in \mathbb{N}$ ,  $e_j$  occurs with probability  $\theta_j^*$ . For every  $n$  and for every  $\varepsilon \in [0, 1]$ , the sample mean  $\bar{X}_n$  satisfies

$$d_{TV}(\bar{X}_n, \theta^*) \leq \frac{1}{\sqrt{n}} \bar{X}_n^{1/2} + 2\varepsilon + 3\sqrt{\frac{\log(2/\delta)}{2n}}$$

with probability at least  $1 - \delta$ , where  $\delta \in (0, 1)$ . We also have

$$\mathbf{E}d_{TV}(\bar{X}_n, \theta^*) \leq \frac{1}{\sqrt{n}} \mathbf{E} \bar{X}_n^{1/2} + 2\varepsilon.$$

## Various models of contamination

**Huber's contamination** ( $\mathcal{M}_n^{HC}(\varepsilon, \theta^*)$ ): The observations  $X_1, \dots, X_n$  are independent and drawn from the mixture distribution

$$P_{\varepsilon, \theta^*, Q} := (1 - \varepsilon)P_{\theta^*} + \varepsilon Q,$$

where  $P_{\theta^*}$  represents the reference distribution parametrized by  $\theta^*$ ,  $Q$  is the distribution of the outliers and  $\varepsilon$  is the fraction of the outliers.

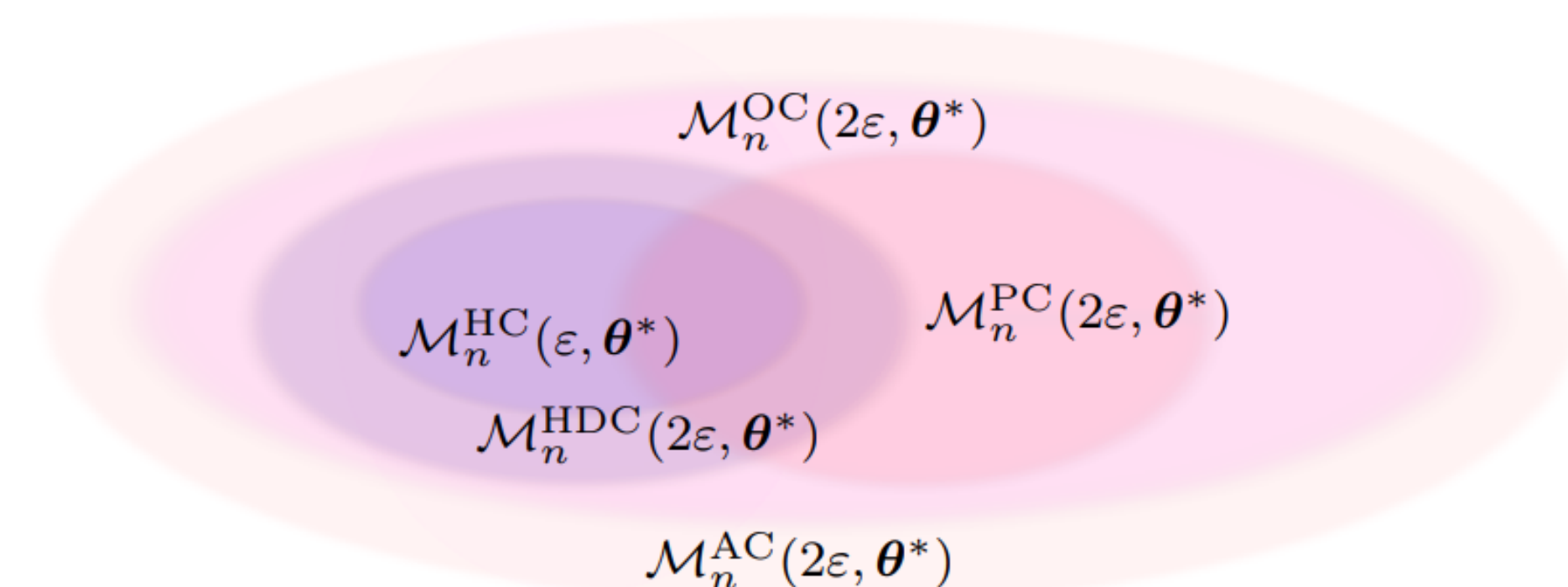
**Huber's deterministic contamination** ( $\mathcal{M}_n^{HDC}(\varepsilon, \theta^*)$ ): The number of outliers is not random, and there are at most  $n\varepsilon$  outliers following the contamination distribution.

**Oblivious contamination** ( $\mathcal{M}_n^{OC}(\varepsilon, \theta^*)$ ): The outliers are not necessarily i.i.d. and the number and the joint distribution of outliers is determined in advance, possibly based on the knowledge of the reference distribution.

**Parameter contamination** ( $\mathcal{M}_n^{PC}(\varepsilon, \theta^*)$ ): The parameters of the distributions of some observations are contaminated. Each outlier  $X_i$  is drawn from a distribution  $Q_i = P_{\theta_i}$  belonging to the same family as the reference distribution, but corresponding to a contaminated parameter  $\theta_i \neq \theta^*$ .

**Adversarial contamination** ( $\mathcal{M}_n^{AC}(\varepsilon, \theta^*)$ ): An adversary sees all the initial observations drawn from the reference distribution and replaces  $n\varepsilon$  of them by arbitrary values.

## Hierarchy between models



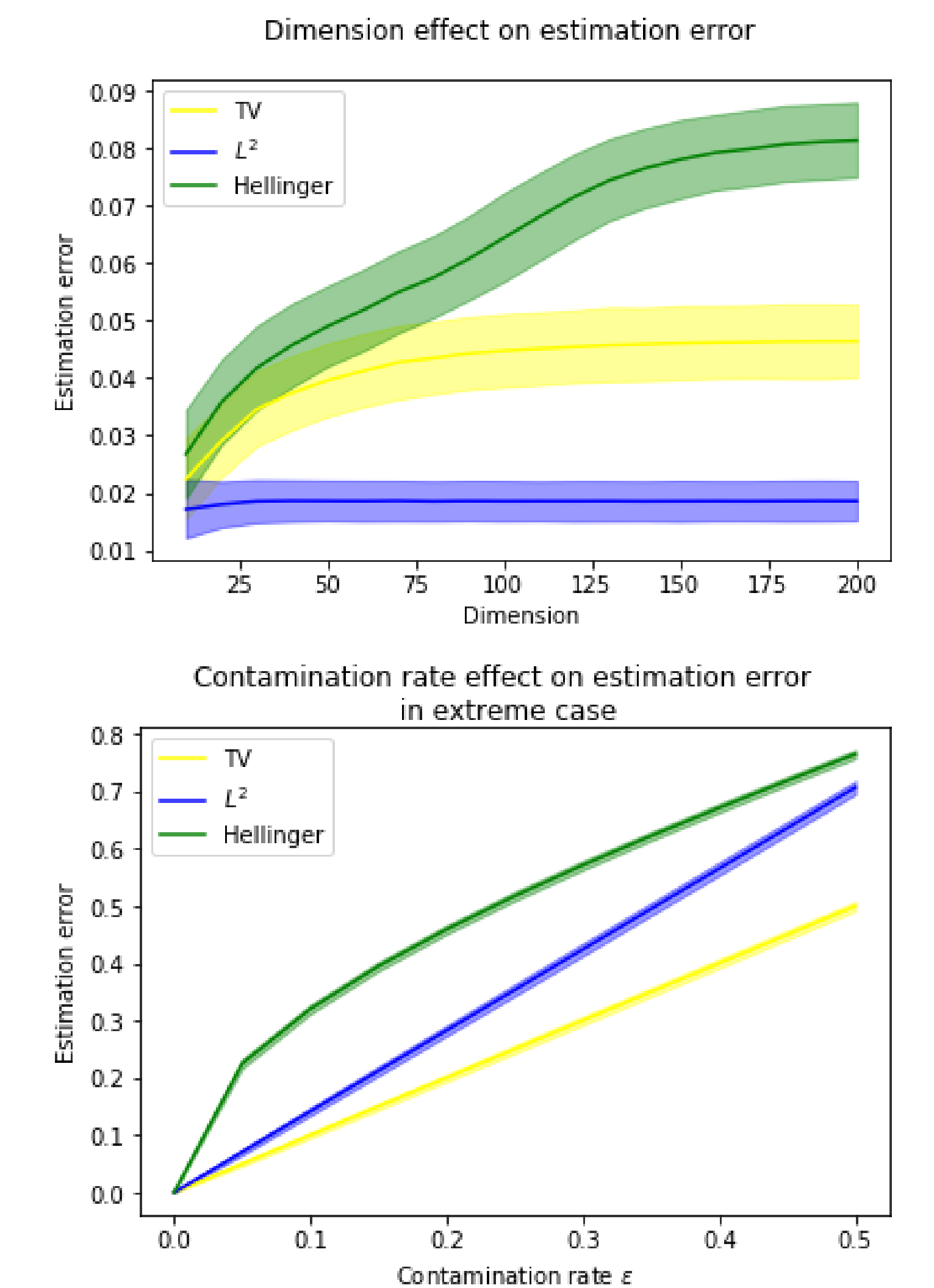
## Confidence regions

**Theorem 3.** Let  $\delta \in (0, 1)$  be the tolerance level. If  $\theta^* \in \Delta_s^{k-1}$ , then under any contamination model, the regions of  $\Delta^{k-1}$  defined by each of the following inequalities

$$\begin{aligned} d_{\mathbb{L}^2}(\bar{X}_n, \theta) &\leq (1/n)^{1/2} + \sqrt{2}\varepsilon + (\log(1/\delta)/n)^{1/2}, \\ d_{TV}(\bar{X}_n, \theta) &\leq (s/n)^{1/2} + 2\varepsilon + (2\log(1/\delta)/n)^{1/2}, \\ d_H(\bar{X}_n, \theta) &\leq \sqrt{5}((s/n)\log(2s/\delta))^{1/2} + \varepsilon^{1/2} + ((1/2n)\log(2/\delta))^{1/2}, \end{aligned}$$

contain  $\theta^*$  with probability at least  $1 - \delta$ .

## Illustration on a numerical example



## Conclusion

The choice of the distance has much stronger impact on the risk rate than the nature of contamination.

## References

Bateni, A.-H. and Dalalyan, A. S. (2020). Confidence regions and minimax rates in outlier-robust estimation on the probability simplex. *Electron. J. Statist.*, 14(2):2653–2677.