

EMPIRICAL RISK MINIMIZATION FOR HEAVY-TAILED LOSSES

BY CHRISTIAN BROWNLEES¹, EMILIE JOLY AND GÁBOR LUGOSI¹

Pompeu Fabra University, HEC Paris–CNRS and Pompeu Fabra University

The purpose of this paper is to discuss empirical risk minimization when the losses are not necessarily bounded and may have a distribution with heavy tails. In such situations, usual empirical averages may fail to provide reliable estimates and empirical risk minimization may provide large excess risk. However, some robust mean estimators proposed in the literature may be used to replace empirical means. In this paper, we investigate empirical risk minimization based on a robust estimate proposed by Catoni. We develop performance bounds based on chaining arguments tailored to Catoni’s mean estimator.

1. Introduction. Heavy-tailed data are commonly encountered in many fields of research (see, e.g., Embrechts, Klüppelberg and Mikosch [14] and Finkenstadt and Rootzén [16]). For instance, in finance, the influential work of Mandelbrot [22] and Fama [15] documented evidence of power-law behavior in asset prices in the early 1960s. When the data have heavy tails, standard statistical procedures typically perform poorly and appropriate robust alternatives are needed to carry out inference effectively. In this paper, we propose a class of robust empirical risk minimization procedures for such data that are based on a robust estimator introduced by Catoni [12].

Empirical risk minimization is one of the basic principles of statistical learning that is routinely applied in a great variety of problems such as regression function estimation, classification and clustering. The general model may be described as follows. Let X be a random variable taking values in some measurable space \mathcal{X} and let \mathcal{F} be a set of nonnegative functions defined on \mathcal{X} . For each $f \in \mathcal{F}$, define the risk $m_f = \mathbb{E}f(X)$ and let $m^* = \inf_{f \in \mathcal{F}} m_f$ denote the optimal risk. In statistical learning, n independent random variables X_1, \dots, X_n are available, all distributed as X , and one aims at finding a function with small risk. To this end, one may define the *empirical risk minimizer*

$$f_{\text{ERM}} = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(X_i),$$

Received June 2014; revised May 2015.

¹Supported by the Spanish Ministry of Science and Technology Grant MTM2012-37195. *MSC2010 subject classifications.* Primary, 62F35; secondary 62F12.

Key words and phrases. Empirical risk minimization, heavy-tailed data, robust regression, robust k -means clustering, Catoni’s estimator.

where, for the simplicity of the discussion and essentially without loss of generality, we implicitly assume that the minimizer exists. If the minimum is achieved by more than one function, one may pick one of them arbitrarily.

REMARK (Loss functions and risks). The main motivation and terminology may be explained by the following general prediction problem in statistical learning. Let the “training data” $(Z_1, Y_1), \dots, (Z_n, Y_n)$ be independent identically distributed pairs of random variables where the Z_i take their values in, say, \mathbb{R}^m and the Y_i are real-valued. In classification problems, the Y_i take discrete values. Given a new observation Z , one is interested in predicting the value of the corresponding response variable Y where the pair (Z, Y) has the same distribution as that of the (Z_i, Y_i) . A predictor is a function $g : \mathbb{R}^m \rightarrow \mathbb{R}$ whose quality is measured with the help of a *loss function* $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$. The *risk* of g is then $\mathbb{E}\ell(g(Z), Y)$. Given a class \mathcal{G} of functions $g : \mathbb{R}^m \rightarrow \mathbb{R}$, empirical risk minimization chooses one that minimizes the *empirical risk* $(1/n) \sum_{i=1}^n \ell(g(Z_i), Y_i)$ over all $g \in \mathcal{G}$. In the simplified notation followed in this paper, X_i corresponds to the pair (Z_i, Y_i) , the function f represents $\ell(g(\cdot), \cdot)$ and m_f substitutes $\mathbb{E}\ell(g(Z), Y)$.

The performance of empirical risk minimization is measured by the *risk* of the selected function,

$$m_{\text{ERM}} = \mathbb{E}[f_{\text{ERM}}(X)|X_1, \dots, X_n].$$

In particular, the main object of interest for this paper is the *excess risk* $m_{\text{ERM}} - m^*$. The performance of empirical risk minimization has been thoroughly studied and well understood using tools of empirical process theory. In particular, the simple observation that

$$m_{\text{ERM}} - m^* \leq 2 \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(X_i) - m_f \right|,$$

allows one to apply the rich theory on the suprema of empirical processes to obtain upper performance bounds. The interested reader is referred to Bartlett and Mendelson [7], Boucheron, Bousquet and Lugosi [9], Koltchinskii [18], Massart [23], Mendelson [26], van de Geer [34] for references and recent results in this area. Essentially all of the theory of empirical minimization assumes either that the functions f are uniformly bounded or that the random variables $f(X)$ have sub-Gaussian tails for all $f \in \mathcal{F}$. For example, when all $f \in \mathcal{F}$ take their values in the interval $[0, 1]$, Dudley’s [13] classical metric-entropy bound, together with standard symmetrization arguments, imply that there exists a universal constant c such that

$$(1) \quad \mathbb{E}m_{\text{ERM}} - m^* \leq \frac{c}{\sqrt{n}} \mathbb{E} \int_0^1 \sqrt{\log N_{\mathbb{X}}(\mathcal{F}, \epsilon)} d\epsilon,$$

where for any $\epsilon > 0$, $N_{\mathbb{X}}(\mathcal{F}, \epsilon)$ is the ϵ -covering number of the class \mathcal{F} under the empirical quadratic distance $d_{\mathbb{X}}(f, g) = (\frac{1}{n} \sum_{i=1}^n (f(X_i) - g(X_i))^2)^{1/2}$, defined as the minimal cardinality N of any set $\{f_1, \dots, f_N\} \subset \mathcal{F}$ such that for all $f \in \mathcal{F}$ there exists an $f_j \in \{f_1, \dots, f_N\}$ with $d_{\mathbb{X}}(f, f_j) \leq \epsilon$. Of course, this is one of the most basic bounds and many important refinements have been established.

A tighter bound may be established by the so-called *generic chaining* method; see Talagrand [32]. Recall the following definition (see, e.g., [32], Definition 1.2.3). Let T be a (pseudo) metric space. An increasing sequence (\mathcal{A}_n) of partitions of T is called *admissible* if for all $n = 0, 1, 2, \dots, \#\mathcal{A}_n \leq 2^{2^n}$. For any $t \in T$, denote by $A_n(t)$ the unique element of \mathcal{A}_n that contains t . Let $\Delta(A)$ denote the diameter of the set $A \subset T$. Define, for $\beta = 1, 2$,

$$\gamma_{\beta}(T, d) = \inf_{\mathcal{A}_n} \sup_{t \in T} \sum_{n \geq 0} 2^{n/\beta} \Delta(A_n(t)),$$

where the infimum is taken over all admissible sequences. Then one has

$$(2) \quad \mathbb{E}m_{\text{ERM}} - m^* \leq \frac{c}{\sqrt{n}} \mathbb{E}\gamma_2(\mathcal{F}, d_{\mathbb{X}}),$$

for some universal constant c . This bound implies (1) as $\gamma_2(\mathcal{F}, d_{\mathbb{X}})$ is bounded by a constant multiple of the entropy integral $\int_0^1 \sqrt{\log N_{\mathbb{X}}(\mathcal{F}, \epsilon)} d\epsilon$ (see, e.g., [32]).

However, when the functions f are no longer uniformly bounded and the random variables $f(X)$ may have a heavy tail, empirical risk minimization may have a much poorer performance. This is simply due to the fact that empirical averages become poor estimates of expected values. Indeed, for heavy-tailed distributions, several estimators of the mean are known to outperform simple empirical averages. It is a natural idea to define a robust version of empirical risk minimization based on minimizing such robust estimators.

In this paper, we focus on an elegant and powerful estimator proposed and analyzed by Catoni [12]. (A version of) Catoni’s estimator may be defined as follows.

Introduce the nondecreasing differentiable *truncation function*

$$(3) \quad \phi(x) = -\mathbb{1}_{\{x < 0\}} \log\left(1 - x + \frac{x^2}{2}\right) + \mathbb{1}_{\{x \geq 0\}} \log\left(1 + x + \frac{x^2}{2}\right).$$

To estimate $m_f = \mathbb{E}f(X)$ for some $f \in \mathcal{F}$, define for all $\mu \in \mathbb{R}$,

$$\hat{r}_f(\mu) = \frac{1}{n\alpha} \sum_{i=1}^n \phi(\alpha(f(X_i) - \mu)),$$

where $\alpha > 0$ is a parameter of the estimator to be specified below. Catoni’s estimator of m_f is defined as the unique value $\hat{\mu}_f$ for which $\hat{r}_f(\hat{\mu}_f) = 0$. [Uniqueness is ensured by the strict monotonicity of $\mu \mapsto \hat{r}_f(\mu)$.] Catoni proves that for any fixed $f \in \mathcal{F}$ and $\delta \in [0, 1]$ such that $n > 2 \log(1/\delta)$, under the only assumption that $\text{Var}(f(X)) \leq v$, the estimator above with

$$\alpha = \sqrt{\frac{2 \log(1/\delta)}{n(v + (2v \log(1/\delta))/(n(1 - (2/n) \log(1/\delta))))}}$$

satisfies that, with probability at least $1 - 2\delta$,

$$(4) \quad |m_f - \widehat{\mu}_f| \leq \sqrt{\frac{2v \log(1/\delta)}{n(1 - (2/n) \log(1/\delta))}}.$$

In other words, the deviations of the estimate exhibit a sub-Gaussian behavior. The price to pay is that the estimator depends both on the upper bound v for the variance and on the prescribed confidence δ via the parameter α .

Catoni also shows that for any $n > 4(1 + \log(1/\delta))$, if $\text{Var}(f(X)) \leq v$, the choice

$$\alpha = \sqrt{\frac{2}{nv}}$$

guarantees that, with probability at least $1 - 2\delta$,

$$(5) \quad |m_f - \widehat{\mu}_f| \leq (1 + \log(1/\delta)) \sqrt{\frac{v}{n}}.$$

Even though we lose the sub-Gaussian tail behavior, the estimator is independent of the required confidence level.

Given such a powerful mean estimator, it is natural to propose an empirical risk minimizer that selects a function from the class \mathcal{F} that minimizes Catoni’s mean estimator. Formally, define

$$\widehat{f} = \arg \min_{f \in \mathcal{F}} \widehat{\mu}_f,$$

where again, for the sake of simplicity we assume that the minimizer exists. (Otherwise one may select an appropriate approximate minimizer and all arguments go through in a trivial way.)

Once again, as a first step of understanding the excess risk $m_{\widehat{f}} - m^*$, we may use the simple bound

$$m_{\widehat{f}} - m^* = (m_{\widehat{f}} - \widehat{\mu}_{\widehat{f}}) + (\widehat{\mu}_{\widehat{f}} - m^*) \leq 2 \sup_{f \in \mathcal{F}} |m_f - \widehat{\mu}_f|.$$

When \mathcal{F} is a finite class of cardinality, say $|\mathcal{F}| = N$, Catoni’s bound may be combined, in a straightforward way, with the union-of-events bound. Indeed, if the estimators $\widehat{\mu}_f$ are defined with parameter

$$\alpha = \sqrt{\frac{2 \log(N/\delta)}{n(v + (2v \log(N/\delta)/(n(1 - (2/n) \log(N/\delta))))}},$$

then, with probability at least $1 - 2\delta$,

$$\sup_{f \in \mathcal{F}} |m_f - \widehat{\mu}_f| \leq \sqrt{\frac{2v \log(N/\delta)}{n(1 - (2/n) \log(N/\delta))}}.$$

Note that this bound requires that $\sup_{f \in \mathcal{F}} \text{Var}(f(X)) \leq v$, that is, the variances are uniformly bounded by a *known* value v . Throughout the paper, we work with this assumption. However, this bound does not take into account the structure of the class \mathcal{F} and it is useless when \mathcal{F} is an infinite class. Our strategy to obtain meaningful bounds is to use *chaining* arguments. However, the extension is nontrivial and the argument becomes more involved. The main results of the paper present performance bounds for empirical minimization of Catoni's estimator based on generic chaining.

REMARK (Median-of-means estimator). Catoni's estimator is not the only one with sub-Gaussian deviations for heavy-tailed distributions. Indeed, the *median-of-means* estimator, proposed by Nemirovsky and Yudin [28] (and also independently by Alon, Matias and Szegedy [2]) has similar performance guarantees as (4). This estimate is obtained by dividing the data in several small blocks, calculating the sample mean within each block, and then taking the median of these means. Hsu and Sabato [17] and Minsker [27] introduce multivariate generalizations of the median-of-means estimator and use it to define and analyze certain statistical learning procedures in the presence of heavy-tailed data. The sub-Gaussian behavior is achieved under various assumptions on the loss function. Such conditions can be avoided here. As an example, we detail applications of our results in Section 4 for three different examples of loss functions. An important advantage of the median-of-means estimate over Catoni's estimate is that the parameter of the estimate (i.e., the number of blocks) only depends on the confidence level δ but not on v and, therefore, no prior upper bound of the variance v is required to compute this estimate. Also, the median-of-means estimate is useful even when the variance is infinite and only a moment of order $1 + \epsilon$ exists for some $\epsilon > 0$ (see Bubeck, Cesa-Bianchi and Lugosi [11]). Lerasle and Oliveira [19] consider empirical minimization of the median-of-means estimator and obtain interesting results in various statistical learning problems. However, to establish metric-entropy bounds for minimization of this mean estimate remains to be a challenge.

The rest of the paper is organized as follows. In Section 2, we state and discuss the main results of the paper. Section 3 is dedicated to the proofs. In Section 4, we describe some applications to regression under the absolute and squared losses and k -means clustering. Finally, in Section 5 we present some simulation results both for regression and k -means clustering. The simulation study gives empirical evidence that the proposed empirical risk minimization procedure improves performance in a significant manner in the presence of heavy-tailed data. Some of the more technical arguments are relegated to the [Appendix](#).

2. Main results. The bounds we establish for the excess risk depend on the geometric structure of the class \mathcal{F} under different distances. The $L_2(P)$ distance

is defined, for $f, f' \in \mathcal{F}$, by

$$d(f, f') = (\mathbb{E}[(f(X) - f'(X))^2])^{1/2}$$

and the L_∞ distance is

$$D(f, f') = \sup_{x \in \mathcal{X}} |f(x) - f'(x)|.$$

We also work with the (random) empirical quadratic distance

$$d_{\mathbb{X}}(f, f') = \left(\frac{1}{n} \sum_{i=1}^n (f(X_i) - f'(X_i))^2 \right)^{1/2}.$$

Denote by f^* a function with minimal expectation

$$f^* = \arg \min_{f \in \mathcal{F}} m_f.$$

Next, we present two results that bound the excess risk $m_{\hat{f}} - m_{f^*}$ of the minimizer \hat{f} of Catoni’s risk estimate in terms of metric properties of the class \mathcal{F} . The first result involves a combination of terms involving the γ_2 and γ_1 functionals under the metrics d and D while the second is in terms of quantiles of γ_2 under the empirical metric $d_{\mathbb{X}}$.

THEOREM 1. *Let \mathcal{F} be a class of nonnegative functions defined on a set \mathcal{X} and let X, X_1, \dots, X_n be i.i.d. random variables taking values in \mathcal{X} . Assume that there exists $v > 0$ such that $\sup_{f \in \mathcal{F}} \text{Var}(f(X)) \leq v$. Let $\delta \in (0, 1/3)$. Suppose that \hat{f} is selected from \mathcal{F} by minimizing Catoni’s mean estimator with parameter α . Then there exists a universal constant L such that, under the condition*

$$6\left(\alpha v + \frac{2 \log(\delta^{-1})}{n\alpha}\right) + L \log(2\delta^{-1}) \left(\frac{\gamma_2(\mathcal{F}, d)}{\sqrt{n}} + \frac{\gamma_1(\mathcal{F}, D)}{n} \right) \leq \frac{1}{\alpha},$$

with probability at least $1 - 3\delta$, the risk of \hat{f} satisfies

$$m_{\hat{f}} - m_{f^*} \leq 6\left(\alpha v + \frac{2 \log(\delta^{-1})}{n\alpha}\right) + L \log(2\delta^{-1}) \left(\frac{\gamma_2(\mathcal{F}, d)}{\sqrt{n}} + \frac{\gamma_1(\mathcal{F}, D)}{n} \right).$$

THEOREM 2. *Assume the hypotheses of Theorem 1. We denote by $\text{diam}_d(\mathcal{F})$ the diameter of the class \mathcal{F} under the distance d . Set Γ_δ such that $\mathbb{P}\{\gamma_2(\mathcal{F}, d_{\mathbb{X}}) > \Gamma_\delta\} \leq \frac{\delta}{8}$. Then there exists a universal constant K such that, under the condition*

$$6\left(\alpha v + \frac{2 \log(\delta^{-1})}{n\alpha}\right) + K \max(\Gamma_\delta, \text{diam}_d(\mathcal{F})) \sqrt{\frac{\log(8/\delta)}{n}} \leq \frac{1}{\alpha},$$

with probability at least $1 - 3\delta$, the risk of \hat{f} satisfies

$$m_{\hat{f}} - m_{f^*} \leq 6\left(\alpha v + \frac{2 \log(\delta^{-1})}{n\alpha}\right) + K \max(\Gamma_\delta, \text{diam}_d(\mathcal{F})) \sqrt{\frac{\log(8/\delta)}{n}}.$$

In both theorems above, the choice of α only influences the term $\alpha v + 2 \log(\delta^{-1})/(\alpha n)$. By taking $\alpha = \sqrt{2 \log(\delta^{-1})/(nv)}$, this term equals

$$2\sqrt{\frac{2v \log(\delta^{-1})}{n}}.$$

For example, in that case, the condition in Theorem 1 reduces to

$$12\sqrt{\frac{2v \log(\delta^{-1})}{n}} + L \log(\delta^{-1}) \left(\frac{\gamma_2(\mathcal{F}, d)}{\sqrt{n}} + \frac{\gamma_1(\mathcal{F}, D)}{n} \right) \leq \sqrt{\frac{nv}{2 \log(\delta^{-1})}}.$$

This holds for sufficiently large values of n . This choice has the disadvantage that the estimator depends on the confidence level (i.e., on the value of δ). By taking $\alpha = \sqrt{2/(nv)}$, independently of δ , one obtains the slightly worse term

$$\sqrt{\frac{2v}{n}}(1 + \log(\delta^{-1})).$$

Observe that the main term in the second part of the bound of Theorem 1 is

$$L \log(\delta^{-1}) \frac{\gamma_2(\mathcal{F}, d)}{\sqrt{n}}$$

which is comparable to the bound (2) obtained under the strong condition of $f(X)$ being uniformly bounded. All other terms are of smaller order. Note that this part of the bound depends on the “weak” distribution-dependent $L_2(P)$ metric d . The quantity $\gamma_1(\mathcal{F}, D) \geq \gamma_2(\mathcal{F}, d)$ also enters the bound of Theorem 1 though only multiplied by $1/n$. The presence of this term requires that \mathcal{F} be bounded in the L_∞ distance D which limits the usefulness of the bound. In Section 4, we illustrate the bounds on two applications to regression and k -means clustering. In these applications, in spite of the presence of heavy tails, the covering numbers under the distance D may be bounded in a meaningful way. Note that no such bound can hold for “ordinary” empirical risk minimization that minimizes the usual empirical means $(1/n) \sum_{i=1}^n f(X_i)$ because of the poor performance of empirical averages in the presence of heavy tails.

The main merit of the bound of Theorem 2 is that it does not require that the class \mathcal{F} has a finite diameter under the supremum norm. Instead, the quantiles of $\gamma_2(\mathcal{F}, d_{\mathbb{X}})$ enter the picture. In Section 4, we show through the example of L_2 regression how these quantiles may be estimated.

3. Proofs. The proofs of Theorems 1 and 2 are based on showing that the excess risk can be bounded as soon as the supremum of the empirical process $\{X_f(\mu) : f \in \mathcal{F}\}$ is bounded for any fixed $\mu \in \mathbb{R}$, where for any $f \in \mathcal{F}$ and $\mu \in \mathbb{R}$, we define $X_f(\mu) = \widehat{r}_f(\mu) - \bar{r}_f(\mu)$ with

$$\bar{r}_f(\mu) = \frac{1}{\alpha} \mathbb{E}[\phi(\alpha(f(X) - \mu))]$$

and

$$\widehat{r}_f(\mu) = \frac{1}{n\alpha} \sum_{i=1}^n \phi(\alpha(f(X_i) - \mu)).$$

The two theorems differ in the way the supremum of this empirical process is bounded.

Let $A_\alpha(\delta) = \alpha v + 2 \log(\delta^{-1})/(n\alpha)$.

Once again, we may assume, essentially without loss of generality, that the minimum exists. In case of multiple minimizers, we may choose one arbitrarily. The main result in [12] states that for any $\delta > 0$ such that $\alpha^2 v + 2 \log(\delta^{-1})/n \leq 1$, with probability at least $1 - 2\delta$,

$$(6) \quad |\widehat{\mu}_{f^*} - m_{f^*}| \leq A_\alpha(\delta).$$

Let $\Omega_{f^*}(\delta)$ be the event on which inequality (6) holds. By definition, $\mathbb{P}\{\Omega_{f^*}(\delta)\} \geq 1 - 2\delta$.

3.1. *A deterministic version of $\widehat{\mu}_f$.* We begin with a variant of the argument of Catoni [12]. It involves a deterministic version $\overline{\mu}_f$ of the estimator defined, for each $f \in \mathcal{F}$, as the unique solution of the equation $\overline{r}_f(\mu) = 0$.

In Lemma 3 below, we show that $\overline{\mu}_f$ is in a small (deterministic) interval centered at m_f . For any $f \in \mathcal{F}$, $\mu \in \mathbb{R}$, and $\varepsilon \geq 0$, define

$$B_f^+(\mu, \varepsilon) = (m_f - \mu) + \frac{\alpha}{2}(m_f - \mu)^2 + \frac{\alpha}{2}v + \varepsilon,$$

$$B_f^-(\mu, \varepsilon) = (m_f - \mu) - \frac{\alpha}{2}(m_f - \mu)^2 - \frac{\alpha}{2}v - \varepsilon$$

and let

$$\mu_f^+(\varepsilon) = m_f + \alpha v + 2\varepsilon, \quad \mu_f^-(\varepsilon) = m_f - \alpha v - 2\varepsilon.$$

As a function of μ , $B_f^+(\mu, \varepsilon)$ is a quadratic polynomial such that $\mu_f^+(\varepsilon)$ is an upper bound of the smallest root of $B_f^+(\mu, \varepsilon)$. Similarly, $\mu_f^-(\varepsilon)$ is a lower bound of the largest root of $B_f^-(\mu, \varepsilon)$. Implicitly, we assumed that these roots always exist. This is not always the case but a simple condition on α guarantees that these roots exists. In particular, $1 - \alpha^2 v - 2\alpha\varepsilon \geq 0$ guarantees that $B_f^+(\mu, \varepsilon) = 0$ and $B_f^-(\mu, \varepsilon) = 0$ have at least one solution. This condition will always be satisfied by our choice of ε and α .

Still following the ideas of [12], the next lemma bounds $\overline{r}_f(\mu)$ by the quadratic polynomials B^+ and B^- . The lemma will help us compare the zero of $\overline{r}_f(\mu)$ to the zeros of these quadratic functions.

LEMMA 3. *For any fixed $f \in \mathcal{F}$ and $\mu \in \mathbb{R}$,*

$$(7) \quad B_f^-(\mu, 0) \leq \overline{r}_f(\mu) \leq B_f^+(\mu, 0),$$

and, therefore, $m_f - \alpha v \leq \bar{\mu}_f \leq m_f + \alpha v$. In particular,

$$B_{\hat{f}}^-(\mu, 0) \leq \bar{r}_{\hat{f}}(\mu) \leq B_{\hat{f}}^+(\mu, 0).$$

For any μ and ε , such that $\bar{r}_{\hat{f}}(\mu) \leq \varepsilon$, if $1 - \alpha^2 v - 2\alpha\varepsilon \geq 0$, then

$$(8) \quad m_{\hat{f}} \leq \mu + \alpha v + 2\varepsilon.$$

PROOF. Writing Y for $\alpha(f(X) - \mu)$ and using the fact that $\phi(x) \leq \log(1 + x + x^2/2)$ for all $x \in \mathbb{R}$,

$$\begin{aligned} \exp(\alpha \bar{r}_f(\mu)) &\leq \exp\left(\mathbb{E}\left[\log\left(1 + Y + \frac{Y^2}{2}\right)\right]\right) \\ &\leq \mathbb{E}\left[1 + Y + \frac{Y^2}{2}\right] \\ &\leq 1 + \alpha(m_f - \mu) + \frac{\alpha^2}{2}[v + (m_f - \mu)^2] \leq \exp(\alpha B_f^+(\mu, 0)). \end{aligned}$$

Thus, we have $\bar{r}_f(\mu) - B_f^+(\mu, 0) \leq 0$ (see Figure 1). Since this last inequality is true for any f , $\sup_f(\bar{r}_f(\mu) - B_f^+(\mu, 0)) \leq 0$ and the second inequality of (7) is proved. The second statement of the lemma may be proved by a similar argument.

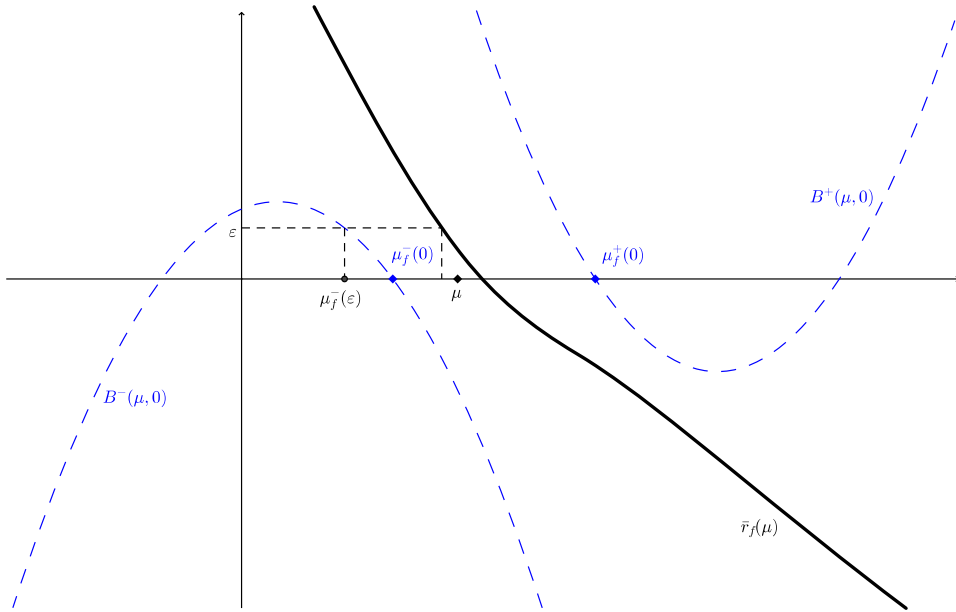


FIG. 1. Representation of $\bar{r}_f(\mu)$ and the quadratic functions $B_f^-(\mu, 0)$ and $B_f^+(\mu, 0)$. $\bar{r}_f(\mu)$ is squeezed between $B_f^-(\mu, 0)$ and $B_f^+(\mu, 0)$. In particular at $\mu_f^+(0)$ [resp., $\mu_f^-(0)$], $\bar{r}_f(\mu)$ is nonpositive (resp., nonnegative). Any μ such that $\bar{r}_f(\mu) \leq \varepsilon$ is above $\mu_f^-(\varepsilon)$.

If $\bar{r}_{\hat{f}}(\mu) \leq \varepsilon$, then $B_{\hat{f}}^-(\mu, 0) \leq \varepsilon$ which is equivalent to $B_{\hat{f}}^-(\mu, \varepsilon) \leq 0$. If $1 - \alpha^2 v - 2\alpha\varepsilon \geq 0$ then a solution of $B_{\hat{f}}^-(\mu, \varepsilon) = 0$ exists and since $\bar{r}_{\hat{f}}(\mu)$ is a nonincreasing function, μ is above the largest of these two solutions. This implies $\mu_{\hat{f}}^-(\varepsilon) \leq \mu$ which gives inequality (8) (see Figure 1). \square

Inequality (8) is the key tool to ensure that the risk $m_{\hat{f}}$ of the minimizer \hat{f} can be upper bounded as soon as $\bar{r}_{\hat{f}}$ is. It remains to find the smallest μ and ε such that $\bar{r}_f(\mu)$ is bounded uniformly on \mathcal{F} .

3.2. *Bounding the excess risk in terms of the supremum of an empirical process.* The key to all proofs is that we link the excess risk to the supremum of the empirical process $X_f(\mu) = \hat{r}_f(\mu) - \bar{r}_f(\mu)$ as f ranges through \mathcal{F} for a suitably chosen value of μ . For fixed $\mu \in \mathbb{R}$ and $\delta \in (0, 1)$, define the $1 - \delta$ quantile of $\sup_{f \in \mathcal{F}} |X_f(\mu) - X_{f^*}(\mu)|$ by $Q(\mu, \delta)$, that is, the infimum of all positive numbers q such that

$$\mathbb{P}\left\{\sup_{f \in \mathcal{F}} |X_f(\mu) - X_{f^*}(\mu)| \leq q\right\} \geq 1 - \delta.$$

First, we need a few simple facts summarized in the next lemma.

LEMMA 4. *Let $\mu_0 = m_{f^*} + A_\alpha(\delta)$. Then on the event $\Omega_{f^*}(\delta)$, the following inequalities hold:*

1. $\hat{r}_{\hat{f}}(\mu_0) \leq 0$;
2. $\bar{r}_{f^*}(\mu_0) \leq 0$;
3. $-\hat{r}_{f^*}(\mu_0) \leq 2A_\alpha(\delta)$.

PROOF. We prove each inequality separately.

1. First, note that on $\Omega_{f^*}(\delta)$ inequality (6) holds, and we have $\hat{\mu}_{\hat{f}} \leq \hat{\mu}_{f^*} \leq \mu_0$. Since $\hat{r}_{\hat{f}}$ is a nonincreasing function of μ , $\hat{r}_{\hat{f}}(\mu_0) \leq \hat{r}_{\hat{f}}(\hat{\mu}_{\hat{f}}) = 0$.

2. By (7), $\bar{\mu}_{f^*} \leq m_{f^*} + \alpha v \leq m_{f^*} + \alpha v + 2 \log(\delta^{-1}) / (n\alpha) = \mu_0$. Since \bar{r}_{f^*} is a nonincreasing function, $\bar{r}_{f^*}(\mu_0) \leq \bar{r}_{f^*}(\bar{\mu}_{f^*}) = 0$.

3. \hat{r}_{f^*} is a 1-Lipschitz function and, therefore,

$$\begin{aligned} |\hat{r}_{f^*}(\mu_0)| &= |\hat{r}_{f^*}(\hat{\mu}_{f^*}) - \hat{r}_{f^*}(\mu_0)| \leq |\hat{\mu}_{f^*} - \mu_0| \\ &\leq |\hat{\mu}_{f^*} - m_{f^*}| + |m_{f^*} - \mu_0| \\ &\leq 2A_\alpha(\delta) \end{aligned}$$

which gives $-\hat{r}_{f^*}(\mu_0) \leq 2A_\alpha(\delta)$. \square

We will use Lemma 3 with μ_0 introduced in Lemma 4. Recall that $\mathbb{P}\{\Omega_{f^*}(\delta)\} \geq 1 - 2\delta$.

With the notation introduced above, we see that with probability at least $1 - \delta$,

$$\begin{aligned} \bar{r}_{\hat{f}}(\mu_0) &\leq \hat{r}_{\hat{f}}(\mu_0) + \bar{r}_{f^*}(\mu_0) - \hat{r}_{f^*}(\mu_0) \\ &\quad + |\bar{r}_{\hat{f}}(\mu_0) - \hat{r}_{\hat{f}}(\mu_0) - \bar{r}_{f^*}(\mu_0) + \hat{r}_{f^*}(\mu_0)| \\ &\leq \hat{r}_{\hat{f}}(\mu_0) + \bar{r}_{f^*}(\mu_0) - \hat{r}_{f^*}(\mu_0) \\ &\quad + \sup_{f \in \mathcal{F}} |\bar{r}_f(\mu_0) - \hat{r}_f(\mu_0) - \bar{r}_{f^*}(\mu_0) + \hat{r}_{f^*}(\mu_0)| \\ &\leq \hat{r}_{\hat{f}}(\mu_0) + \bar{r}_{f^*}(\mu_0) - \hat{r}_{f^*}(\mu_0) + Q(\mu_0, \delta). \end{aligned}$$

This inequality, together with Lemma 4, implies that, with probability at least $1 - 3\delta$,

$$\bar{r}_{\hat{f}}(\mu_0) \leq 2A_\alpha(\delta) + Q(\mu_0, \delta).$$

Now using Lemma 3 with $\varepsilon = 2A_\alpha(\delta) + Q(\mu_0, \delta)$ and under the condition $1 - \alpha^2 v - 4\alpha A_\alpha(\delta) - 2\alpha Q(\mu_0, \delta) \geq 0$, we have

$$\begin{aligned} (9) \quad m_{\hat{f}} - m_{f^*} &\leq \alpha v + 5A_\alpha(\delta) + 2Q(\mu_0, \delta) \\ &\leq 6\left(\alpha v + \frac{2 \log(\delta^{-1})}{n\alpha}\right) + 2Q(\mu_0, \delta), \end{aligned}$$

with probability at least $1 - 3\delta$. The condition $1 - \alpha^2 v - 4\alpha A_\alpha(\delta) - 2\alpha Q(\mu_0, \delta) \geq 0$ is satisfied whenever

$$6\left(\alpha v + \frac{2 \log(\delta^{-1})}{n\alpha}\right) + 2Q(\mu_0, \delta) \leq \frac{1}{\alpha}$$

holds.

3.3. Bounding the supremum of the empirical process. Theorems 1 and 2 both follow from (9) by two different ways of bounding the quantile $Q(\mu, \delta)$ of $\sup_{f \in \mathcal{F}} |X_f(\mu) - X_{f^*}(\mu)|$. Here, we present these two inequalities. Both of them use basic results of “generic chaining”; see Talagrand [32]. Theorem 1 follows from (9) and the next inequality.

PROPOSITION 5. *Let $\mu \in \mathbb{R}$ and $\alpha > 0$. There exists a universal constant L such that for any $\delta \in (0, 1)$,*

$$Q(\mu, \delta) \leq L \log(2\delta^{-1}) \left(\frac{\gamma_2(\mathcal{F}, d)}{\sqrt{n}} + \frac{\gamma_1(\mathcal{F}, D)}{n} \right).$$

The proof is an immediate consequence of Theorem 12 and (14) in the Appendix and the following lemma.

LEMMA 6. For any $\mu \in \mathbb{R}$, $\alpha > 0$, $f, f' \in \mathcal{F}$, and $t > 0$,

$$\mathbb{P}\{|X_f(\mu) - X_{f'}(\mu)| > t\} \leq 2 \exp\left(-\frac{nt^2}{2(d(f, f')^2 + (2D(f, f')t/(3)))}\right),$$

where the distances d, D are defined at the beginning of Section 2.

PROOF. Observe that $n(X_f(\mu) - X_{f'}(\mu))$ is the sum of the independent zero-mean random variables

$$C_i(f, f') = \frac{1}{\alpha}\phi(\alpha(f(X_i) - \mu)) - \frac{1}{\alpha}\phi(\alpha(f'(X_i) - \mu)) - \left[\frac{1}{\alpha}\mathbb{E}[\phi(\alpha(f(X) - \mu))] - \frac{1}{\alpha}\mathbb{E}[\phi(\alpha(f'(X) - \mu))]\right].$$

Note that since the truncation function ϕ is 1-Lipschitz, we have $C_i(f, f') \leq 2D(f, f')$. Also,

$$\sum_{i=1}^n \mathbb{E}[C_i(f, f')^2] \leq \sum_{i=1}^n \mathbb{E}[\left((f(X_i) - \mu) - (f'(X_i) - \mu)\right)^2] = nd(f, f')^2.$$

The lemma follows from Bernstein’s inequality [see, e.g., [10], equation (2.10)]. □

Similarly, Theorem 2 is implied by (9) and the following. Recall the notation of Theorem 2.

THEOREM 7. Let $\mu \in \mathbb{R}$, $\alpha > 0$, and $\delta \in (0, 1/3)$. There exists a universal constant K such that

$$Q(\mu, \delta) \leq K \max(\Gamma_\delta, \text{diam}_d(\mathcal{F})) \sqrt{\frac{\log(8/\delta)}{n}}.$$

PROOF. Assume $\Gamma_\delta \geq \text{diam}_d(\mathcal{F})$. The proof is based on a standard symmetrization argument. Let (X'_1, \dots, X'_n) be independent copies of (X_1, \dots, X_n) and define

$$Z_i(f) = \frac{1}{n\alpha}\phi(\alpha(f(X_i) - \mu)) - \frac{1}{n\alpha}\phi(\alpha(f(X'_i) - \mu)).$$

Introduce also independent Rademacher random variables $(\varepsilon_1, \dots, \varepsilon_n)$. For any $f \in \mathcal{F}$, denote by $Z(f) = \sum_{i=1}^n \varepsilon_i Z_i(f)$. Then by Hoeffding’s inequality, for all $f, g \in \mathcal{F}$ and for every $t > 0$,

$$(10) \quad \mathbb{P}_{(\varepsilon_1, \dots, \varepsilon_n)}\{|Z(f) - Z(g)| > t\} \leq 2 \exp\left(-\frac{t^2}{2d_{\mathbb{X}, \mathbb{X}'}(f, g)^2}\right),$$

where $\mathbb{P}_{(\varepsilon_1, \dots, \varepsilon_n)}$ denotes probability with respect to the Rademacher variables only (i.e., conditional on the X_i and X'_i) and $d_{\mathbb{X}, \mathbb{X}'}(f, g) = \sqrt{\sum_{i=1}^n (Z_i(f) - Z_i(g))^2}$ is a random distance. Using (16) in the Appendix with distance $d_{\mathbb{X}, \mathbb{X}'}$ and (10), we get that, for all $\lambda > 0$,

$$(11) \quad \mathbb{E}_{(\varepsilon_1, \dots, \varepsilon_n)} \left[\exp \left(\lambda \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \varepsilon_i [Z_i(f) - Z_i(f^*)] \right| \right) \right] \leq 2 \exp(\lambda^2 L^2 \gamma_2(\mathcal{F}, d_{\mathbb{X}, \mathbb{X}'})^2 / 4),$$

where L is a universal constant from Proposition 14. Observe that since $x \mapsto \phi(x)$ is Lipschitz with constant 1,

$$\begin{aligned} d_{\mathbb{X}, \mathbb{X}'}(f, g) &= \left(\frac{1}{n^2 \alpha^2} \sum_{i=1}^n (\phi(\alpha(f(X_i) - \mu)) - \phi(\alpha(f(X'_i) - \mu)) - \phi(\alpha(g(X_i) - \mu)) + \phi(\alpha(g(X'_i) - \mu)))^2 \right)^{1/2} \\ &\leq \frac{1}{\sqrt{n}} \left(\frac{1}{n} \sum_{i=1}^n (f(X_i) - g(X_i))^2 \right)^{1/2} + \frac{1}{\sqrt{n}} \left(\frac{1}{n} \sum_{i=1}^n (f(X'_i) - g(X'_i))^2 \right)^{1/2}. \end{aligned}$$

This implies

$$\gamma_2(\mathcal{F}, d_{\mathbb{X}, \mathbb{X}'}) \leq \frac{1}{\sqrt{n}} (\gamma_2(\mathcal{F}, d_{\mathbb{X}}) + \gamma_2(\mathcal{F}, d_{\mathbb{X}'})).$$

Combining this with (11), we obtain

$$\begin{aligned} &\mathbb{P} \left\{ \sup_{f \in \mathcal{F}} |Z(f) - Z(f^*)| \geq t \right\} \\ &\leq \mathbb{P} \left\{ \sup_{f \in \mathcal{F}} |Z(f) - Z(f^*)| \geq t \mid \gamma_2(\mathcal{F}, d_{\mathbb{X}}) \leq \Gamma_\delta \text{ and } \gamma_2(\mathcal{F}, d_{\mathbb{X}'}) \leq \Gamma_\delta \right\} \\ &\quad + 2\mathbb{P} \{ \gamma_2(\mathcal{F}, d_{\mathbb{X}}) > \Gamma_\delta \} \\ &\leq \mathbb{E}_{\mathbb{X}, \mathbb{X}'} \left[\mathbb{E}_{(\varepsilon_1, \dots, \varepsilon_n)} \left[e^{\lambda \sup_{f \in \mathcal{F}} \left| \sum_{i=1}^n \varepsilon_i [Z_i(f) - Z_i(f^*)] \right|} \mid \gamma_2(\mathcal{F}, d_{\mathbb{X}}) \leq \Gamma_\delta \text{ and } \gamma_2(\mathcal{F}, d_{\mathbb{X}'}) \leq \Gamma_\delta \right] e^{-\lambda t} \right. \\ &\quad \left. + \frac{\delta}{4} \quad (\text{by the definition of } \Gamma_\delta) \right] \\ &\leq 2 \exp \left(\frac{\lambda^2 L^2}{n} \Gamma_\delta^2 - \lambda t \right) + \frac{\delta}{4}. \end{aligned}$$

Optimization in λ with $t = 2L\Gamma_\delta\sqrt{\log(8/\delta)/n}$ gives

$$\mathbb{P}\left\{\sup_{f \in \mathcal{F}} |Z(f) - Z(f^*)| \geq t\right\} \leq \frac{\delta}{2}.$$

A standard symmetrization inequality of tail probabilities of empirical processes (see, e.g., [34], Lemma 3.3) guarantees that

$$\mathbb{P}\left\{\sup_{f \in \mathcal{F}} |X_f(\mu) - X_{f^*}(\mu)| \geq 2t\right\} \leq 2\mathbb{P}\left\{\sup_{f \in \mathcal{F}} |Z(f) - Z(f^*)| \geq t\right\}$$

as long as for any $f \in \mathcal{F}$, $\mathbb{P}\{|X_f(\mu) - X_{f^*}(\mu)| \geq t\} \leq \frac{1}{2}$. Recall that $X_f(\mu) - X_{f^*}(\mu)$ is a zero-mean random variable. Then by Chebyshev’s inequality, it suffices to have $t \geq \sqrt{2} \text{diam}_d(\mathcal{F})/\sqrt{n}$. Indeed,

$$\begin{aligned} & \frac{\text{Var}(X_f(\mu) - X_{f^*}(\mu))}{t^2} \\ & \leq \frac{\text{Var}((1/\alpha)\phi(\alpha(f(X) - \mu)) - (1/\alpha)\phi(\alpha(f^*(X) - \mu)))}{nt^2} \\ & \leq \frac{\mathbb{E}[(f(X) - f^*(X))^2]}{nt^2} \\ & \leq \frac{\text{diam}_d(\mathcal{F})^2}{nt^2}. \end{aligned}$$

Without loss of generality, we can assume $L \geq 1$. Since for any choice of $\delta < \frac{1}{3}$, $\sqrt{\log(\frac{8}{\delta})} > \sqrt{2}$ we have $L\Gamma_\delta\sqrt{\log(\frac{8}{\delta})} \geq \text{diam}_d(\mathcal{F})\sqrt{2}$. Thus,

$$\mathbb{P}\left\{\sup_{f \in \mathcal{F}} |X_f(\mu) - X_{f^*}(\mu)| \geq 2L\Gamma_\delta\sqrt{\frac{\log(8/\delta)}{n}}\right\} \leq \delta$$

as desired. Now, if $\Gamma_\delta < \text{diam}_d(\mathcal{F})$, $\mathbb{P}\{\gamma_2(\mathcal{F}, d_{\mathbb{X}}) > \text{diam}_d(\mathcal{F})\} \leq \frac{\delta}{8}$ and the same argument holds for $\text{diam}_d(\mathcal{F})$ instead of Γ_δ . This completes the proof. \square

4. Applications. In this section, we describe two applications of Theorems 1 and 2 to simple statistical learning problems. The first is a regression estimation problem in which we distinguish between L_1 and L_2 risks and the second is k -means clustering.

4.1. *Empirical risk minimization for regression.*

4.1.1. *L_1 regression.* Let $(Z_1, Y_1), \dots, (Z_n, Y_n)$ be independent identically distributed random variables taking values in $\mathcal{Z} \times \mathbb{R}$ where \mathcal{Z} is a bounded subset of (say) \mathbb{R}^m . Suppose \mathcal{G} is a class of functions $\mathcal{Z} \rightarrow \mathbb{R}$ bounded in the L_∞ norm, that is, $\sup_{g \in \mathcal{G}} \sup_{z \in \mathcal{Z}} |g(z)| < \infty$. We denote by Δ the diameter of \mathcal{G} under the

distance induced by this norm. First, we consider the setup when the *risk* of each $g \in \mathcal{G}$ is defined by the L_1 loss

$$R(g) = \mathbb{E}|g(Z) - Y|,$$

where the pair (Z, Y) has the same distribution of the (Z_i, Y_i) and is independent of them. Let $g^* = \arg \min_{g \in \mathcal{G}} R(g)$ be a minimizer of the risk (which, without loss of generality, is assumed to exist). The statistical learning problem we consider here consists of choosing a function \hat{g} from the class \mathcal{G} that has a risk $R(\hat{g})$ not much larger than $R(g^*)$.

The standard procedure is to pick \hat{g} by minimizing the empirical risk $(1/n) \sum_{i=1}^n |g(Z_i) - Y_i|$ over $g \in \mathcal{G}$. However, if the response variable Y is unbounded and may have a heavy tail, ordinary empirical risk minimization may fail to provide a good predictor of Y as the empirical risk is an unreliable estimate of the true risk.

Here, we propose choosing \hat{g} by minimizing Catoni’s estimate. To this end, we only need to assume that the second moment of Y is bounded by a known constant. More precisely, assume that $\mathbb{E}Y^2 \leq \sigma^2$ for some $\sigma > 0$. Then $\sup_{g \in \mathcal{G}} \text{Var}(|g(Z) - Y|) \leq 2\sigma^2 + 2 \sup_{g \in \mathcal{G}} \sup_{z \in \mathcal{Z}} |g(z)|^2 \stackrel{\text{def}}{=} v$ is a known and finite constant.

Now for all $g \in \mathcal{G}$ and $\mu \in \mathbb{R}$, define

$$\hat{r}_g(\mu) = \frac{1}{n\alpha} \sum_{i=1}^n \phi(\alpha(|g(X_i) - Y_i| - \mu)),$$

where ϕ is the truncation function defined in (3). Define $\hat{R}(g)$ as the unique value for which $\hat{r}_g(\hat{R}(g)) = 0$. The empirical risk minimizer based on Catoni’s risk estimate is then

$$\hat{g} = \arg \min_{g \in \mathcal{G}} \hat{R}(g).$$

By Theorem 1, the performance of \hat{g} may be bounded in terms of covering numbers of the class of functions $\mathcal{F} = \{f(z, y) = |g(z) - y| : g \in \mathcal{G}\}$ based on the distance

$$D(f, f') = \sup_{z \in \mathcal{Z}, y \in \mathbb{R}} ||g(z) - y| - |g'(z) - y|| \leq \sup_{z \in \mathcal{Z}} |g(z) - g'(z)|.$$

Thus, the covering numbers of \mathcal{F} under the distance D may be bounded in terms of the covering numbers of \mathcal{G} under the L_∞ distance. Denoting by $N_d(A, \epsilon)$ the ϵ -covering number of a set A under the metric d , we obtain the following.

COROLLARY 8. *Consider the setup described above. We assume $\int_0^\Delta \log N_\infty(\mathcal{G}, \epsilon) d\epsilon < \infty$. Let $n \in \mathbb{N}$, $\delta \in (0, 1/3)$ and $\alpha = \sqrt{2 \log(\delta^{-1})/(nv)}$.*

There exists an integer N_0 and a universal constant C such that, for all $n \geq N_0$, with probability at least $1 - 3\delta$,

$$R(\hat{g}) - R(g^*) \leq 12\sqrt{\frac{2v \log(\delta^{-1})}{n}} + C \log(2\delta^{-1}) \left(\frac{1}{\sqrt{n}} \int_0^\Delta \sqrt{\log N_d(\mathcal{G}, \epsilon)} d\epsilon + O\left(\frac{1}{n}\right) \right).$$

PROOF. Clearly, if two distances d_1 and d_2 satisfy $d_1 \leq d_2$, then $\gamma_1(\mathcal{F}, d_1) \leq \gamma_1(\mathcal{F}, d_2)$. Thus, $\gamma_1(\mathcal{F}, D) \leq \gamma_1(\mathcal{G}, \|\cdot\|_\infty) \leq L \int_0^\Delta \log N_\infty(\mathcal{G}, \epsilon) d\epsilon < \infty$ [see (15)] and $\gamma_1(\mathcal{F}, D)/n = O(1/n)$. The condition

$$12\sqrt{\frac{2v \log(\delta^{-1})}{n}} + C \log(2\delta^{-1}) \left(\frac{1}{\sqrt{n}} \int_0^\Delta \sqrt{\log N_d(\mathcal{G}, \epsilon)} d\epsilon + O\left(\frac{1}{n}\right) \right) \leq \sqrt{\frac{nv}{2 \log(\delta^{-1})}}$$

is satisfied for sufficiently large n . Apply Theorem 1. \square

Note that the bound essentially has the same form as (1) but to apply (1) it is crucial that the response variable Y is bounded or at least has sub-Gaussian tails. We get this under the weak assumption that Y has a bounded second moment (with a known upper bound). The price we pay is that covering numbers under the distance $d_{\mathbb{X}}$ are now replaced by covering numbers under the supremum norm.

4.1.2. *L₂ regression.* Here, we consider the same setup as in Section 4.1.1 but now the risk is measured by the L_2 loss. The *risk* of each $g \in \mathcal{G}$ is defined by the L_2 loss

$$R(g) = \mathbb{E}(g(Z) - Y)^2.$$

Note that Theorem 1 is useless here as the difference $|R(g) - R(g')|$ is not bounded by the L_∞ distance of g and g' anymore and the covering numbers of \mathcal{F} under the metric D are infinite. However, Theorem 2 gives meaningful bounds. Let $g^* = \arg \min_{g \in \mathcal{G}} R(g)$ and again we choose \hat{g} by minimizing Catoni’s estimate.

Here, we need to assume that $\mathbb{E}Y^4 \leq \sigma^2$ for some $\sigma > 0$. Then $\sup_{g \in \mathcal{G}} \text{Var}((g(Z) - Y)^2) \leq 8\sigma^2 + 8 \sup_{g \in \mathcal{G}} \sup_{z \in \mathcal{Z}} |g(z)|^4 \stackrel{\text{def}}{=} v$ is a known and finite constant.

By Theorem 2, the performance of \hat{g} may be bounded in terms of covering numbers of the class of functions $\mathcal{F} = \{f(z, y) = (g(z) - y)^2 : g \in \mathcal{G}\}$ based on the distance

$$d_{\mathbb{X}}(f, f') = \left(\frac{1}{n} \sum_{i=1}^n ((g(Z_i) - Y_i)^2 - (g'(Z_i) - Y_i)^2)^2 \right)^{1/2}.$$

Note that

$$\begin{aligned} |(g(Z_i) - Y_i)^2 - (g'(Z_i) - Y_i)^2| &= |g(Z_i) - g'(Z_i)| |2Y_i - g(Z_i) - g'(Z_i)| \\ &\leq 2|g(Z_i) - g'(Z_i)| (|Y_i| + \Delta) \\ &\leq 2d_\infty(g, g') (|Y_i| + \Delta), \end{aligned}$$

and, therefore,

$$\begin{aligned} d_{\mathbb{X}}(f, f') &\leq 2d_\infty(g, g') \sqrt{\frac{1}{n} \sum_{i=1}^n (|Y_i| + \Delta)^2} \\ &\leq 2\sqrt{2}d_\infty(g, g') \sqrt{\Delta^2 + \frac{1}{n} \sum_{i=1}^n Y_i^2}. \end{aligned}$$

By Chebyshev’s inequality,

$$\mathbb{P}\left\{ \frac{1}{n} \sum_{i=1}^n Y_i^2 - \mathbb{E}[Y^2] > t \right\} \leq \frac{\text{Var}(Y^2)}{nt^2} \leq \frac{\sigma^2}{nt^2}$$

thus $\frac{1}{n} \sum_{i=1}^n Y_i^2 > \mathbb{E}[Y^2] + \sqrt{8\sigma^2/(n\delta)}$ with probability at most $\delta/8$ and

$$d_{\mathbb{X}}(f, f') > 2\sqrt{2}d_\infty(g, g') \sqrt{\Delta^2 + \mathbb{E}[Y^2] + \sqrt{\frac{8\sigma^2}{n\delta}}}$$

occurs with a probability bounded by $\frac{\delta}{8}$. Recall again that for two distances d_1 and d_2 such that $d_1 \leq cd_2$ one has $\gamma_2(\mathcal{G}, d_1) \leq c\gamma_2(\mathcal{G}, d_2)$. Then Theorem 2 applies with

$$\Gamma_\delta = 2\sqrt{2} \sqrt{\Delta^2 + \mathbb{E}[Y^2] + \sqrt{\frac{8\sigma^2}{n\delta}}} \gamma_2(\mathcal{G}, d_\infty)$$

and $\Gamma_\delta \geq \Delta \geq \text{diam}_d(\mathcal{F})$.

COROLLARY 9. *Consider the setup described above. Let $n \in \mathbb{N}$, $\delta \in (0, 1/3)$ and $\alpha = \sqrt{2 \log(\delta^{-1})/(nv)}$. There exists an integer N_0 and a universal constant C such that, for all $n \geq N_0$, with probability at least $1 - 3\delta$,*

$$\begin{aligned} R(\widehat{g}) - R(g^*) &\leq 12\sqrt{\frac{2v \log(\delta^{-1})}{n}} \\ &\quad + C\sqrt{\log\left(\frac{8}{\delta}\right)} \sqrt{\frac{\Delta^2 + \mathbb{E}[Y^2] + 8\sigma^2/(n\delta)}{n}} \int_0^\Delta \sqrt{\log N_\infty(\mathcal{G}, \epsilon)} d\epsilon. \end{aligned}$$

PROOF. Apply Theorem 2 and note that the condition holds for sufficiently large n . \square

The bound of the corollary essentially matches the best rates of convergence one can get even in the case of bounded regression under such general conditions. For special cases, such as linear regression, better bounds may be proven for other methods; see Audibert and Catoni [5], Hsu and Sabato [17] and Minsker [27].

4.2. *k-means clustering under heavy-tailed distribution.* In *k-means clustering*—or *vector quantization*—one wishes to represent a distribution by a finite number of points. Formally, let X be a random vector taking values in \mathbb{R}^m and let P denote the distribution of X . Let $k \geq 2$ be a positive integer that we fix for the rest of the section. A clustering scheme is given by a set of k cluster centers $C = \{y_1, \dots, y_k\} \subset \mathbb{R}^m$ and a *quantizer* $q : \mathbb{R}^m \rightarrow C$. Given a *distortion measure* $\ell : \mathbb{R}^m \times \mathbb{R}^m \rightarrow [0, \infty)$, one wishes to find C and q such that the expected distortion

$$D_k(P, q) = \mathbb{E}\ell(X, q(X))$$

is as small as possible. The minimization problem is meaningful whenever $\mathbb{E}\ell(X, 0) < \infty$ which we assume throughout. Typical distortion measures are of the form $\ell(x, y) = \|x - y\|^\alpha$ where $\|\cdot\|$ is a norm on \mathbb{R}^m and $\alpha > 0$ (typically α equals 1 or 2). Here, for concreteness and simplicity, we assume that ℓ is the Euclidean distance $\ell(x, y) = \|x - y\|$ though the results may be generalized in a straightforward manner to other norms. In a way equivalent to the arguments of Section 4.1.2, the results may be generalized to the case of the quadratic distortion $\ell(x, y) = \|x - y\|^2$. In order to avoid repetition of arguments, the details are omitted.

It is not difficult to see that if $\mathbb{E}\|X\| < \infty$, then there exists a (not necessarily unique) quantizer q^* that is optimal, that is, q^* is such that for all clustering schemes q ,

$$D_k(P, q) \geq D_k(P, q^*) \stackrel{\text{def}}{=} D_k^*(P).$$

It is also clear that q^* is a *nearest neighbor quantizer*, that is,

$$\|x - q^*(x)\| = \min_{y_i \in C} \|x - y_i\|.$$

Thus, nearest neighbor quantizers are determined by their cluster centers $C = \{y_1, \dots, y_k\}$. In fact, for all quantizers with a particular set C of cluster centers, the corresponding nearest neighbor quantizer has minimal distortion and, therefore, it suffices to restrict our attention to nearest neighbor quantizers.

In the problem of empirical quantizer design, one is given an i.i.d. sample X_1, \dots, X_n drawn from the distribution P and one's aim is to find a quantizer q_n whose distortion

$$D_k(P, q_n) = \mathbb{E}[\|X - q_n(X)\| | X_1, \dots, X_n]$$

is as close to $D_k^*(P)$ as possible. A natural strategy is to choose a quantizer—or equivalently, a set C of cluster centers—by minimizing the *empirical distortion*

$$D_k(P_n, q) = \frac{1}{n} \sum_{i=1}^n \|X_i - q(X_i)\| = \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, k} \|X_i - y_j\|,$$

where P_n denotes the standard empirical distribution based on X_1, \dots, X_n . If $\mathbb{E}\|X\| < \infty$, then the empirically optimal quantizer asymptotically minimizes the distortion. More precisely, if q_n denotes the empirically optimal quantizer [i.e., $q_n = \arg \min_q D_k(P_n, q)$], then

$$\lim_{n \rightarrow \infty} D_k(P, q_n) = D_k^*(P) \quad \text{with probability 1;}$$

see Pollard [29, 31] and Abaya and Wise [1] (see also Linder [21]). The rate of convergence of $D_k(P, q_n)$ to $D_k^*(P)$ has drawn considerable attention; see, for example, Pollard [30], Bartlett, Linder and Lugosi [6], Antos [3], Antos, Györfi and Györfy [4], Biau, Devroye and Lugosi [8], Maurer and Pontil [25] and Levrard [20]. Such rates are typically studied under the assumption that X is almost surely bounded. Under such assumptions, one can show that

$$\mathbb{E}D_k(P, q_n) - D_k^*(P) \leq C(P, k, m)n^{-1/2},$$

where the constant $C(P, k, m)$ depends on $\text{ess sup } \|X\|$, k , and the dimension m . The value of the constant has mostly been investigated in the case of quadratic loss $\ell(x, y) = \|x - y\|^2$ but most proofs may be modified for the case studied here. For the quadratic loss, one may take $C(P, k, m)$ as a constant multiple of $B^2 \min(\sqrt{k^{1-2/m}m}, k)$ where $B = \text{ess sup } \|X\|$.

However, little is known about the finite-sample performance of empirically designed quantizers under possibly heavy-tailed distributions. In fact, there is no hope to extend the results cited above for distributions with finite second moment simply because empirical averages are poor estimators of means under such general conditions.

In the recent paper of Telgarsky and Dasgupta [33], bounds on the excess risk under conditions on higher moments have been developed. They prove a bound of $\mathcal{O}(n^{-1/2+2/p})$ for the excess distortion where p is the number of moments of $\|X\|$ that are assumed to be finite. Here, we show that there exists an empirical quantizer \hat{q}_n whose excess distortion $D_k(P, \hat{q}_n) - D_k^*(P)$ is of the order of $n^{-1/2}$ (with high probability) under the only assumption that $\mathbb{E}[\|X\|^2]$ is finite. This may be achieved by choosing a quantizer that minimizes Catoni's estimate of the distortion.

The proposed empirical quantizer uses two parameters that depend on the (unknown) distribution of X . For simplicity, we assume that upper bounds for these two parameters are available. (Otherwise either one may try to estimate them or, as the sample size grows, use increasing values for these parameters. The details go beyond the scope of this paper.)

One of these parameters is the second moment $\text{Var}(X) = \mathbb{E}[\|X - \mathbb{E}[X]\|^2]$ and let V be an upper bound. The other parameter $\rho > 0$ is an upper bound for the norm of the possible cluster centers. The next lemma offers an estimate.

LEMMA 10 (Linder [21]). *Let $2 \leq j \leq k$ be the unique integer such that $D_k^* = \dots = D_j^* < D_{j-1}^*$ and define $\varepsilon = (D_{j-1}^* - D_j^*)/2$. Let (y_1, \dots, y_j) be a set of cluster centers such that the distortion of the corresponding quantizer is less than $D_j^* + \varepsilon$. Let $B_r = \{x : \|x\| \leq r\}$ denote the closed ball of radius $r > 0$ centered at the origin. If $\rho > 0$, is such that:*

- $\frac{\rho}{10} P(B_{\rho/10}) > 2\mathbb{E}\|X\|,$
- $P(B_{2\rho/5}) > 1 - \frac{\varepsilon^2}{4\mathbb{E}[\|X\|^2]},$

then for all $1 \leq i \leq k, \|y_i\| \leq \rho.$

Now we are prepared to describe the proposed empirical quantizer. Let \mathcal{C}_ρ be the set of all collections $C = \{y_1, \dots, y_k\} \in (\mathbb{R}^m)^k$ of cluster centers with $\|y_j\| \leq \rho$ for all $j = 1, \dots, k$. For each $C \in \mathcal{C}_\rho$, denote by q_C the corresponding quantizer. Now for all $C \in \mathcal{C}_\rho$, we may calculate Catoni’s mean estimator of the distortion $D(P, q_C) = \mathbb{E}\|X - q_C(X)\| = \mathbb{E} \min_{j=1, \dots, k} \|X_i - y_j\|$ defined as the unique value $\mu \in \mathbb{R}$ for which

$$\frac{1}{n\alpha} \sum_{i=1}^n \phi\left(\alpha \left(\min_{j=1, \dots, k} \|X_i - y_j\| - \mu\right)\right) = 0,$$

where we use the parameter value $\alpha = \sqrt{2/(nkV)}$. Denote this estimator by $\widehat{D}(P_n, q_C)$ and let \widehat{q}_n be any quantizer minimizing the estimated distortion. An easy compactness argument shows that such a minimizer exists.

The main result of this section is the following bound for the distortion of the chosen quantizer.

THEOREM 11. *Assume that $\text{Var}(X) \leq V < \infty$ and $n \geq m$. Then, with probability at least $1 - \delta,$*

$$D_k(P, \widehat{q}_n) - D_k(P, q^*) \leq C \left(\log \frac{1}{\delta}\right) \left(\sqrt{\frac{Vk}{n}} + \sqrt{\frac{mk}{n}}\right) + O\left(\frac{1}{n}\right),$$

where the constant C only depends on $\rho.$

PROOF. The result follows from Theorem 1. All we need to check is that $\text{Var}(\min_{j=1, \dots, k} \|X - y_j\|)$ is bounded by kV and estimate the covering numbers of the class of functions

$$\mathcal{F}_\rho = \left\{ f_C(x) = \min_{y \in C} \|x - y\| : C \in \mathcal{C}_\rho \right\}.$$

The variance bound follows simply by the fact that for all $C \in \mathcal{C}$,

$$\begin{aligned} \text{Var}\left(\min_{j=1,\dots,k} \|X - y_j\|\right) &\leq \sum_{i=1}^k \text{Var}(\|X - y_i\|) \\ &\leq \sum_{i=1}^k \mathbb{E}[\|X - \mathbb{E}X\|^2] + \|\mathbb{E}X - y_i\|^2 - \mathbb{E}[\|X - y_i\|]^2 \\ &\leq kV. \end{aligned}$$

In order to use the bound of Theorem 1, we need to bound the covering numbers of the class \mathcal{F}_ρ under both metrics d and D . We begin with the metric

$$D(f_C, f_{C'}) = \sup_{x \in \mathbb{R}^m} |f_C(x) - f_{C'}(x)|.$$

The notation $B_z(\epsilon, d)$ refers to the ball under the metric d of radius ϵ centered at z . Let Z be a subset of B_ρ such that

$$\mathcal{B}_{B_\rho} := \{B_z(\epsilon, \|\cdot\|) : z \in Z\}$$

is a covering of the set B_ρ by balls of radius ϵ under the Euclidean norm. Let $C \in \mathcal{C}_\rho$ and associate to any $y_i \in C$ one of the centers in Z such that $\|y_i - z_i\| \leq \epsilon$. If there is more than one possible choice for z_i , we pick one of them arbitrarily. We denote by $q_{C'}$ the nearest neighbor quantizer with codebook $C' = (z_i)_i$. Finally, let $S_i = q_{C'}^{-1}(z_i)$. Now clearly, $\forall i, \forall x \in S_i$

$$\begin{aligned} f_C(x) - f_{C'}(x) &= \min_{1 \leq j \leq k} \|x - y_j\| - \min_{1 \leq j \leq k} \|x - z_j\| \\ &= \min_{1 \leq j \leq k} \|x - y_j\| - \|x - z_i\| \\ &\leq \|x - y_i\| - \|x - z_i\| \leq \epsilon \end{aligned}$$

and similarly, $f_{C'}(x) - f_C(x) \leq \epsilon$. Then $f_C \in B_{f_{C'}}(\epsilon, D)$ and

$$\mathcal{B}_{\mathcal{F}_\rho} := \{B_{f_C}(\epsilon, D) : C \in \mathcal{Z}^k\}$$

is a covering of \mathcal{F}_ρ . Since Z can be taken such that $|Z| = N_{\|\cdot\|}(B_\rho, \epsilon)$ we obtain

$$N_d(\mathcal{F}_\rho, \epsilon) \leq N_D(\mathcal{F}_\rho, \epsilon) \leq N_{\|\cdot\|}(B_\rho, \epsilon)^k.$$

By standard estimates on the covering numbers of the ball B_ρ by balls of size ϵ under the Euclidean metric,

$$N_{\|\cdot\|}(B_\rho, \epsilon) \leq \left(\frac{4\rho}{\epsilon}\right)^m$$

(see, e.g., Matousek [24]). In other words, there exists a universal constant L and constants C_ρ and C'_ρ that depends only on ρ such that

$$\gamma_2(\mathcal{F}_\rho, d) \leq L \int_0^{2\rho} \sqrt{\log N_d(\mathcal{F}_\rho, \epsilon)} d\epsilon \leq C_\rho \sqrt{km},$$

and

$$\gamma_1(\mathcal{F}_\rho, D) \leq L \int_0^{2\rho} \log N_D(\mathcal{F}_\rho, \epsilon) d\epsilon \leq C'_\rho km.$$

Theorem 1 may now be applied to the class \mathcal{F}_ρ . \square

5. Simulation study. In this closing section, we present the results of two simulation exercises that assess the performance of the estimators developed in this work.

5.1. *L₂ regression.* The first application is an L_2 regression exercise. Data are simulated from a linear model with heavy-tailed errors and the L_2 regression procedure based on Catoni’s risk minimizer introduced in Section 4.1.2 is used for estimation. The procedure is benchmarked against regular (“vanilla”) L_2 regression based on the minimization of the empirical L_2 loss.

The simulation exercise is designed as follows. We simulate $(Z_1, Y_1), (Z_2, Y_2), \dots, (Z_n, Y_n)$ i.i.d. pairs of random variables in $\mathbb{R}^5 \times \mathbb{R}$. The vector Z_i of explanatory variables is drawn from a multivariate normal distribution with zero mean, unit variance and correlation matrix equal to an equi-correlation matrix with correlation $\rho = 0.9$. The response variable Y_i is generated as

$$Y_i = Z_i^T \theta^* + \epsilon_i,$$

where the parameter vector θ^* is set to $(0.25, -0.25, 0.50, 0.70, -0.75)$ and ϵ_i is a zero mean error term. The error term ϵ_i is drawn from a Pareto distribution with tail parameter β and is appropriately recentered in order to have zero mean. As it is well known, the tail parameter β determines which moments of the Pareto random variable are finite. More specifically, the moment of order k exists only if $k < \beta$. The focus is on finding the value of θ which minimizes the L_2 risk

$$\mathbb{E}|Y - Z_i^T \theta|^2.$$

The parameter θ is estimated using the Catoni and the vanilla L_2 regressions. Let $\widehat{R}_C(\theta)$ denote the solution of the equation

$$\widehat{r}_\theta(\mu) = \frac{1}{n\alpha} \sum_{i=1}^n \phi(\alpha(|Y_i - Z_i^T \theta|^2 - \mu)) = 0.$$

Then the Catoni L_2 regression estimator is defined as

$$\widehat{\theta}_{nC} = \arg \min_{\theta} \widehat{R}_C(\theta).$$

The vanilla L_2 regression estimator is defined as the minimizer of the empirical L_2 loss,

$$\widehat{\theta}_{nV} = \arg \min_{\theta} \widehat{R}_V(\theta) = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n |Y_i - Z_i^T \theta|^2,$$

which is the classical least squares estimator. The precision of each estimator is measured by their excess risk

$$R(\widehat{\theta}_{nC}) - R(\theta^*) = \mathbb{E}|Y - Z^T \widehat{\theta}_{nC}|^2 - \mathbb{E}|Y - Z^T \theta^*|^2,$$

$$R(\widehat{\theta}_{nV}) - R(\theta^*) = \mathbb{E}|Y - Z^T \widehat{\theta}_{nV}|^2 - \mathbb{E}|Y - Z^T \theta^*|^2.$$

We estimate excess risk by simulation. For each replication of the simulation exercise, we estimate the risk of the estimators and the optimal risk using sample averages based on an i.i.d. sample $(Z'_1, Y'_1), \dots, (Z'_m, Y'_m)$ that is independent of the one used for estimation, that is,

$$\begin{aligned} \widetilde{R}(\widehat{\theta}_{nC}) &= \frac{1}{m} \sum_{i=1}^m |Y'_i - Z_i'^T \widehat{\theta}_{nC}|^2, \\ \widetilde{R}(\widehat{\theta}_{nV}) &= \frac{1}{m} \sum_{i=1}^m |Y'_i - Z_i'^T \widehat{\theta}_{nV}|^2, \\ \widetilde{R}(\theta^*) &= \frac{1}{m} \sum_{i=1}^m |Y'_i - Z_i'^T \theta^*|^2. \end{aligned} \tag{12}$$

The simulation experiment is replicated for different values of the Pareto tail parameter β ranging from 2.01 to 6.01 and different values of the sample size n , ranging from 50 to 1000. For each combination of the tail parameter β and sample size n , the experiment is replicated 1000 times.

Figure 2 displays the Monte Carlo estimate of the excess risk of the Catoni and benchmark regression estimators as functions of the tail parameter β when the sample size n is equal to 500. The left panel shows the level of the excess risks $R(\widehat{\theta}_{nC}) - R(\theta^*)$ and $R(\widehat{\theta}_{nV}) - R(\theta^*)$ as a function of β and the right one shows the percentage improvement of the excess risk of the Catoni procedure over the benchmark calculated as $(R(\widehat{\theta}_{nV}) - R(\widehat{\theta}_{nC})) / (R(\widehat{\theta}_{nC}) - R(\theta^*))$. When the tails are

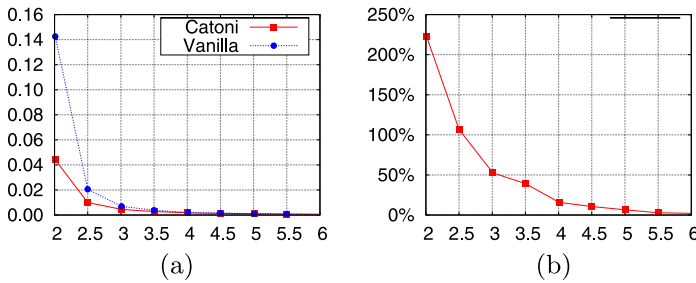


FIG. 2. L_2 regression parameter estimation. The figure plots the excess risk of the Catoni and vanilla L_2 regression parameter estimators (a) and the percentage improvement of the Catoni procedure relative to the vanilla (b) as a function of the tail parameter β for a sample size n equal to 500.

TABLE 1
Relative performance of the Catoni L_2 parameter estimator

β	$n = 50$	$n = 100$	$n = 250$	$n = 500$	$n = 750$	$n = 1000$
2.01	3872.10	440.50	171.30	222.70	218.20	142.80
2.50	169.20	158.70	151.50	106.70	91.70	57.40
3.01	137.60	178.00	89.00	52.50	62.70	63.50
3.50	54.40	20.90	41.30	39.20	38.10	33.50
4.01	30.20	44.40	25.50	15.70	16.30	15.90
4.50	16.50	12.10	11.30	10.60	6.90	13.70
5.01	10.20	7.80	10.20	6.40	5.70	3.10
5.50	6.00	14.80	3.90	2.90	2.10	2.20
6.01	3.90	1.90	2.70	2.10	1.90	1.40

The table reports the percentage improvement of the excess risk of the Catoni L_2 regression estimator relative to the vanilla procedure for different values of the tail parameter β and sample size n .

not excessively heavy (high values of β) the difference between the procedures is small. As the tails become heavier (small values of β), the risks of both procedures increase. Importantly, the Catoni estimator becomes progressively more efficient as the tails become heavier and becomes significantly more efficient when the tail parameter is close to 2. Detailed results for different values of n are reported in Table 1. Overall, the Catoni L_2 regression estimator never performs worse than the benchmark, and it is substantially better when the tails of the data are heavy.

5.2. *k-means.* In the second experiment, we carry out a k -means clustering exercise. Data are simulated from a heavy-tailed mixture distribution and then cluster centers are chosen by minimizing Catoni’s estimate of the L_2 distortion. The performance of the algorithm is benchmarked against the (“vanilla”) k -means algorithm procedure where the distortion is estimated by the standard empirical average.

The simulation exercise is designed as follows. An i.i.d. sample of random vectors X_1, \dots, X_n in \mathbb{R}^2 is drawn from a four-component mixture distribution with equal weights. The means of the mixture components are $(5, 5)$, $(-5, 5)$, $(-5, -5)$ and $(5, -5)$. Each component of the mixture is made up of two appropriately centered independent draws from a Pareto distribution with tail parameter β . The cluster centers obtained by the k -means algorithm based on Catoni and the vanilla k -means algorithm are denoted, respectively, by \hat{q}_{nC} and \hat{q}_{nV} . (Since finding the empirically optimal cluster centers is computationally prohibitive, we use the well-known iterative optimization procedure “ k -means” for the vanilla version and a similar variant for the Catoni scheme.) Analogously to the previous exercise, we summarize the performance of the clustering procedures using the excess risk of the algorithms, that is,

$$D_k(P, \hat{q}_{nC}) - D_k(P, q^*), \quad D_k(P, \hat{q}_{nV}) - D_k(P, q^*),$$

where q^* denotes the means of the mixture components. We estimate excess risk by simulation. We compute the distortion of the quantizers using an i.i.d. sample X'_1, \dots, X'_m of vectors that is independent of the ones used for estimation, that is,

$$\begin{aligned}
 D_k(P_m, \hat{q}_{nC}) &= \frac{1}{m} \sum_{i=1}^m \min_{j=1, \dots, k} \|X'_i - \hat{q}_{nC}(X'_i)\|^2, \\
 D_k(P_m, \hat{q}_{nV}) &= \frac{1}{m} \sum_{i=1}^m \min_{j=1, \dots, k} \|X'_i - \hat{q}_{nV}(X'_i)\|^2, \\
 D_k(P_m, q^*) &= \frac{1}{m} \sum_{i=1}^m \min_{j=1, \dots, k} \|X'_i - q^*(X'_i)\|^2.
 \end{aligned}
 \tag{13}$$

The experiment is replicated for different values of the tail parameter β ranging from 2.01 to 6.01 and different values of the sample size n ranging from 50 to 1000. For each combination of tail parameter β and sample size n the experiment is replicated 1000 times.

Figure 3 displays the Monte Carlo estimate of excess risk of the Catoni and benchmark estimators as a function of tail parameter β for $n = 500$. The left panel shows the estimated excess risk while the right panel shows the percentage improvement of the excess risk of the Catoni procedure, calculated as $(D_k(P, \hat{q}_{nV}) - D_k(P, \hat{q}_{nC})) / (D_k(P, \hat{q}_{nC}) - D_k(P, q^*))$.

The overall results are analogous to the ones of the L_2 regression application. When the tails of the mixture are not excessively heavy (high values of β) the difference in the procedures is small. As the tails become heavier (small values of β), the risk of both procedures increases, but the Catoni algorithm becomes progressively more efficient. The percentage gains for the Catoni procedure are substantial when the tail parameter is smaller than 4. Table 2 reports detailed results for different values of n . As in the L_2 regression simulation study, the Catoni k -means algorithm never performs worse than the benchmark and it is substantially better when the tails of the mixture are heavy.

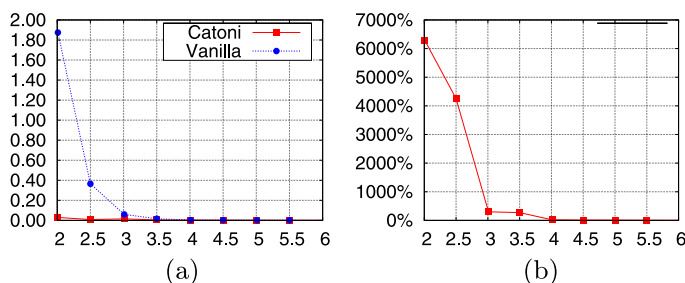


FIG. 3. k -means quantizer estimation. The figure plots the excess risk of the Catoni and vanilla k -means quantizer estimator (a) and the percentage improvement of the Catoni procedure relative to the vanilla (b) as a function of the tail parameter β for a sample size n equal to 500.

TABLE 2
Relative performance of the Catoni k -means quantizer estimator

β	$n = 50$	$n = 100$	$n = 250$	$n = 500$	$n = 750$	$n = 1000$
2.01	823.80	2180.40	3511.60	6278.90	7858.70	10,684.60
2.50	404.50	1007.40	2959.80	4255.40	6828.60	9093.60
3.01	301.10	312.20	286.80	298.60	813.60	1560.20
3.50	129.60	188.60	213.30	271.40	448.60	410.00
4.01	73.80	30.90	26.80	20.30	18.20	13.10
4.50	27.60	22.90	16.50	11.70	9.50	10.10
5.01	16.40	10.80	11.60	8.70	6.00	7.20
5.50	9.00	6.80	9.20	5.00	4.10	4.00
6.01	3.50	4.70	5.00	2.70	3.20	3.10

The table reports the improvement of the Catoni k -means quantizer estimator relative to the vanilla procedure for different values of the tail parameter β and sample size n .

APPENDIX

A.1. A chaining theorem. The following result is a version of standard bounds based on “generic chaining”; see Talagrand [32]. We include the proof for completeness.

Recall that if ψ is a nonnegative increasing convex function defined on \mathbb{R}_+ with $\psi(0) = 0$, then the Orlicz norm of a random variable X is defined by

$$\|X\|_\psi = \inf \left\{ c > 0 : \mathbb{E} \left[\psi \left(\frac{|X|}{c} \right) \right] \leq 1 \right\}.$$

We consider Orlicz norms defined by

$$\psi_1(x) = \exp(x) - 1 \quad \text{and} \quad \psi_2(x) = \exp(x^2) - 1.$$

For further information on Orlicz norms, see [35], Chapter 2.2. First, $\|X\|_{\psi_1} \leq \|X\|_{\psi_2} \sqrt{\log(2)}$ holds. Also note that, by Markov’s inequality, $\|X\|_{\psi_1} \leq c$ implies that $\mathbb{P}\{|X| > t\} \leq 2e^{-t/c}$ and similarly, if $\|X\|_{\psi_2} \leq c$, then $\mathbb{P}\{|X| > t\} \leq 2e^{-t^2/c^2}$. Then

$$(14) \quad \begin{aligned} X &\leq \|X\|_{\psi_1} \log(2\delta^{-1}) && \text{with probability at least } 1 - \delta, \\ X &\leq \|X\|_{\psi_2} \sqrt{\log(2\delta^{-1})} && \text{with probability at least } 1 - \delta. \end{aligned}$$

Recall the following definition (see, e.g., [32], Definition 1.2.3). Let T be a (pseudo) metric space. An increasing sequence (\mathcal{A}_n) of partitions of T is called *admissible* if for all $n = 0, 1, 2, \dots, \#\mathcal{A}_n \leq 2^{2^n}$. For any $t \in T$, denote by $A_n(t)$ the unique element of \mathcal{A}_n that contains t . Let $\Delta(A)$ denote the diameter of the set $A \subset T$. Define, for $\beta = 1, 2$,

$$\gamma_\beta(T, d) = \inf_{\mathcal{A}_n} \sup_{t \in T} \sum_{n \geq 0} 2^{n/\beta} \Delta(A_n(t)),$$

where the infimum is taken over all admissible sequences. First of all, we know from [32], equation (1.18), that there exists a universal constant L such that

$$(15) \quad \gamma_\beta(T, d) \leq L \int_0^{\text{diam}_d(T)} (\log N_d(T, \varepsilon))^{1/\beta} d\varepsilon.$$

THEOREM 12. *Let $(X_t)_{t \in T}$ be a stochastic process indexed by a set T on which two (pseudo) metrics, d_1 and d_2 , are defined such that T is bounded with respect to both metrics. Assume that for any $s, t \in T$ and for all $x > 0$,*

$$\mathbb{P}\{|X_s - X_t| > x\} \leq 2 \exp\left(-\frac{1}{2} \frac{x^2}{d_2(s, t)^2 + d_1(s, t)x}\right).$$

Then for all $t \in T$,

$$\left\| \sup_{s \in T} |X_s - X_t| \right\|_{\psi_1} \leq L(\gamma_1(T, d_1) + \gamma_2(T, d_2))$$

with $L \leq 384 \log(2)$.

The proof of Theorem 12 uses the following lemma.

LEMMA 13 ([35], Lemma 2.2.10). *Let $a, b > 0$ and assume that the random variables X_1, \dots, X_m satisfy, for all $x > 0$,*

$$\mathbb{P}\{|X_i| > x\} \leq 2 \exp\left(-\frac{1}{2} \frac{x^2}{b + ax}\right).$$

Then

$$\left\| \max_{1 \leq i \leq m} X_i \right\|_{\psi_1} \leq 48(a \log(1 + m) + \sqrt{b} \sqrt{\log(1 + m)}).$$

PROOF OF THEOREM 12. Consider an admissible sequence $(\mathcal{B}_n)_{n \geq 0}$ such that for all $t \in T$,

$$\sum_{n \geq 0} 2^n \Delta_1(\mathcal{B}_n(t)) \leq 2\gamma_1(T, d_1)$$

and an admissible sequence $(\mathcal{C}_n)_{n \geq 0}$ such that for all $t \in T$,

$$\sum_{n \geq 0} 2^{n/2} \Delta_2(\mathcal{C}_n(t)) \leq 2\gamma_2(T, d_2).$$

Now we may define an admissible sequence by intersection of the elements of $(\mathcal{B}_{n-1})_{n \geq 1}$ and $(\mathcal{C}_{n-1})_{n \geq 1}$: set $\mathcal{A}_0 = \{T\}$ and let

$$\mathcal{A}_n = \{B \cap C : B \in \mathcal{B}_{n-1} \text{ and } C \in \mathcal{C}_{n-1}\}.$$

$(\mathcal{A}_n)_{n \geq 0}$ is an admissible sequence because each \mathcal{A}_n is increasing and contains at most $(2^{2^{n-1}})^2 = 2^{2^n}$ sets. Define a sequence of finite sets $T_0 = \{t\} \subset T_1 \subset \dots \subset T$

such that T_n contains a single point in each set of \mathcal{A}_n . For any $s \in T$, denote by $\pi_n(s)$ the unique element of T_n in $A_n(s)$. Now for any $s \in T_{k+1}$, we write

$$X_s - X_t = \sum_{k=0}^{\infty} (X_{\pi_{k+1}(s)} - X_{\pi_k(s)}).$$

Then, using the fact that $\|\cdot\|_{\psi_1}$ is a norm and Lemma 13,

$$\begin{aligned} & \left\| \sup_{s \in T} |X_s - X_t| \right\|_{\psi_1} \\ & \leq \sum_{k=0}^{\infty} \left\| \max_{s \in T_{k+1}} |X_{\pi_{k+1}(s)} - X_{\pi_k(s)}| \right\|_{\psi_1} \\ & \leq 48 \sum_{k=0}^{\infty} (d_1(\pi_{k+1}(s), \pi_k(s)) \log(1 + 2^{2^{k+1}}) \\ & \quad + d_2(\pi_{k+1}(s), \pi_k(s)) \sqrt{\log(1 + 2^{2^{k+1}})}). \end{aligned}$$

Since $(\mathcal{A}_n)_{n \geq 0}$ is an increasing sequence, $\pi_{k+1}(s)$ and $\pi_k(s)$ are both in $A_k(s)$. By construction, $A_k(s) \subset B_k(s)$ and, therefore, $d_1(\pi_{k+1}(s), \pi_k(s)) \leq \Delta_1(B_k(s))$. Similarly, we have $d_2(\pi_{k+1}(s), \pi_k(s)) \leq \Delta_2(C_k(s))$. Using $\log(1 + 2^{2^{k+1}}) \leq 4 \log(2)2^k$, we get

$$\begin{aligned} \left\| \max_{s \in T} |X_s - X_t| \right\|_{\psi_1} & \leq 192 \log(2) \left[\sum_{k=0}^{\infty} 2^k \Delta_1(B_k(s)) + \sum_{k=0}^{\infty} 2^{k/2} \Delta_2(C_k(s)) \right] \\ & \leq 384 \log(2) [\gamma_1(T, d_1) + \gamma_2(T, d_2)]. \quad \square \end{aligned}$$

PROPOSITION 14. Assume that for any $s, t \in T$ and for all $x > 0$,

$$\mathbb{P}\{|X_s - X_t| > x\} \leq 2 \exp\left(-\frac{x^2}{2d_2(s, t)^2}\right).$$

Then for all $t \in T$,

$$\left\| \sup_{s \in T} |X_s - X_t| \right\|_{\psi_2} \leq L \gamma_2(T, d_2),$$

where L is a universal constant.

The proof of Proposition 14 is similar to the proof of Theorem 12. One merely needs to replace Lemma 13 by Lemma 2.2.2 in [35] and proceed identically. The details are omitted.

We may use Proposition 14 to bound the moment generating function of $\sup_{s \in T} |X_s - X_t|$ as follows. Set $S = \sup_{s \in T} |X_s - X_t|$. Then using $ab \leq (a^2 + b^2)/2$, we have, for every $\lambda > 0$,

$$\exp(\lambda S) \leq \exp(S^2 / \|S\|_{\psi_2}^2 + \lambda^2 \|S\|_{\psi_2}^2 / 4),$$

and, therefore,

$$(16) \quad \mathbb{E} \left[\exp \left(\lambda \sup_{s \in T} |X_s - X_t| \right) \right] \leq 2 \exp(\lambda^2 L^2 \gamma_2(T, d_2)^2 / 4).$$

Acknowledgments. We thank the referees for their thorough reading and insightful comments that helped greatly improve the manuscript.

REFERENCES

- [1] ABAYA, E. F. and WISE, G. L. (1984). Convergence of vector quantizers with applications to optimal quantization. *SIAM J. Appl. Math.* **44** 183–189. [MR0730008](#)
- [2] ALON, N., MATIAS, Y. and SZEGEDY, M. (1999). The space complexity of approximating the frequency moments. *J. Comput. System Sci.* **58** 137–147. [MR1688610](#)
- [3] ANTOS, A. (2005). Improved minimax bounds on the test and training distortion of empirically designed vector quantizers. *IEEE Trans. Inform. Theory* **51** 4022–4032. [MR2239018](#)
- [4] ANTOS, A., GYÖRFI, L. and GYÖRGY, A. (2005). Individual convergence rates in empirical vector quantizer design. *IEEE Trans. Inform. Theory* **51** 4013–4022. [MR2239017](#)
- [5] AUDIBERT, J.-Y. and CATONI, O. (2011). Robust linear least squares regression. *Ann. Statist.* **39** 2766–2794. [MR2906886](#)
- [6] BARTLETT, P. L., LINDER, T. and LUGOSI, G. (1998). The minimax distortion redundancy in empirical quantizer design. *IEEE Trans. Inform. Theory* **44** 1802–1813. [MR1664098](#)
- [7] BARTLETT, P. L. and MENDELSON, S. (2006). Empirical minimization. *Probab. Theory Related Fields* **135** 311–334. [MR2240689](#)
- [8] BIAU, G., DEVROYE, L. and LUGOSI, G. (2008). On the performance of clustering in Hilbert spaces. *IEEE Trans. Inform. Theory* **54** 781–790. [MR2444554](#)
- [9] BOUCHERON, S., BOUSQUET, O. and LUGOSI, G. (2005). Theory of classification: A survey of some recent advances. *ESAIM Probab. Stat.* **9** 323–375. [MR2182250](#)
- [10] BOUCHERON, S., LUGOSI, G. and MASSART, P. (2013). *Concentration Inequalities: A Nonasymptotic Theory of Independence, with a Foreword by Michel Ledoux*. Oxford Univ. Press, Oxford. [MR3185193](#)
- [11] BUBECK, S., CESA-BIANCHI, N. and LUGOSI, G. (2013). Bandits with heavy tail. *IEEE Trans. Inform. Theory* **59** 7711–7717. [MR3124669](#)
- [12] CATONI, O. (2012). Challenging the empirical mean and empirical variance: A deviation study. *Ann. Inst. Henri Poincaré Probab. Stat.* **48** 1148–1185. [MR3052407](#)
- [13] DUDLEY, R. M. (1978). Central limit theorems for empirical measures. *Ann. Probab.* **6** 899–929. [MR0512411](#)
- [14] EMBRECHTS, P., KLÜPPELBERG, C. and MIKOSCH, T. (1997). *Modelling Extremal Events for Insurance and Finance*. Springer, Berlin. [MR1458613](#)
- [15] FAMA, E. F. (1963). Mandelbrot and the stable Paretian hypothesis. *The Journal of Business* **36** 420–429.
- [16] FINKENSTADT, B. and ROOTZÉN, H. (2003). *Extreme Values in Finance, Telecommunications and the Environment*. Chapman & Hall, New York.
- [17] HSU, D. and SABATO, S. (2013). Approximate loss minimization with heavy tails. Preprint. Available at [arXiv:1307.1827](#).
- [18] KOLTCHINSKII, V. (2006). Local Rademacher complexities and oracle inequalities in risk minimization. *Ann. Statist.* **34** 2593–2656. [MR2329442](#)
- [19] LERASLE, M. and OLIVEIRA, R. I. (2012). Robust empirical mean estimators. Manuscript.

- [20] LEVRARD, C. (2013). Fast rates for empirical vector quantization. *Electron. J. Stat.* **7** 1716–1746. [MR3080408](#)
- [21] LINDER, T. (2002). Learning-theoretic methods in vector quantization. In *Principles of Non-parametric Learning (Udine, 2001)* (L. Györfi, ed.). *CISM Courses and Lectures* **434** 163–210. Springer, Vienna. [MR1987659](#)
- [22] MANDELBROT, B. (1963). The variation of certain speculative prices. *The Journal of Business* **36** 394–419.
- [23] MASSART, P. (2007). *Concentration Inequalities and Model Selection*. *Lecture Notes in Math.* **1896**. Springer, Berlin. [MR2319879](#)
- [24] MATOUŠEK, J. (2002). *Lectures on Discrete Geometry*. *Graduate Texts in Mathematics* **212**. Springer, New York. [MR1899299](#)
- [25] MAURER, A. and PONTIL, M. (2010). K -dimensional coding schemes in Hilbert spaces. *IEEE Trans. Inform. Theory* **56** 5839–5846. [MR2808936](#)
- [26] MENDELSON, S. (2014). Learning without concentration. Preprint. Available at [arXiv:1401.0304](#).
- [27] MINSKER, S. (2015). Geometric median and robust estimation in Banach spaces. *Bernoulli* **21** 2308–2335. [MR3378468](#)
- [28] NEMIROVSKY, A. S. and YUDIN, D. B. (1983). *Problem Complexity and Method Efficiency in Optimization*. Wiley, New York. [MR0702836](#)
- [29] POLLARD, D. (1981). Strong consistency of k -means clustering. *Ann. Statist.* **9** 135–140. [MR0600539](#)
- [30] POLLARD, D. (1982). A central limit theorem for k -means clustering. *Ann. Probab.* **10** 919–926. [MR0672292](#)
- [31] POLLARD, D. (1982). Quantization and the method of k -means. *IEEE Trans. Inform. Theory* **28** 199–205. [MR0651814](#)
- [32] TALAGRAND, M. (2005). *The Generic Chaining*. Springer, Berlin. [MR2133757](#)
- [33] TELGARSKY, M. and DASGUPTA, S. (2013). Moment-based uniform deviation bounds for k -means and friends. Preprint. Available at [arXiv:1311.1903](#).
- [34] VAN DE GEER, S. A. (2000). *Applications of Empirical Process Theory*. Cambridge Univ. Press, Cambridge. [MR1739079](#)
- [35] VAN DER VAART, A. W. and WELLNER, J. A. (1996). *Weak Convergence and Empirical Processes*. Springer, New York. [MR1385671](#)

C. BROWNLEES
 DEPARTMENT OF ECONOMICS AND BUSINESS
 POMPEU FABRA UNIVERSITY
 RAMON TRIAS FARGAS 25-27
 08005 BARCELONA
 SPAIN
 E-MAIL: christian.brownlees@upf.edu

E. JOLY
 GREGHEC
 HEC PARIS–CNRS
 1 RUE DE LA LIBÉRATION 78350
 JOUY-EN-JOSAS
 FRANCE
 E-MAIL: emilien.joly@ens.fr

G. LUGOSI
 ICREA AND DEPARTMENT OF ECONOMICS AND BUSINESS
 POMPEU FABRA UNIVERSITY
 RAMON TRIAS FARGAS 25-27
 08005 BARCELONA
 SPAIN
 E-MAIL: gabor.lugosi@upf.edu