

SUB-GAUSSIAN ESTIMATORS OF THE MEAN OF A RANDOM VECTOR

BY GÁBOR LUGOSI^{*,†,‡,1} AND SHAHAR MENDELSON^{§,¶,2}

ICREA^{}, Pompeu Fabra University[†], Barcelona GSE[‡],
 Technion, I.I.T.[§] and The Australian National University[¶]*

We study the problem of estimating the mean of a random vector X given a sample of N independent, identically distributed points. We introduce a new estimator that achieves a purely sub-Gaussian performance under the only condition that the second moment of X exists. The estimator is based on a novel concept of a multivariate median.

1. Introduction. In this paper, we study the problem of estimating the mean of a random vector X taking values in \mathbb{R}^d . Denoting the mean by $\mu = \mathbb{E}X$, we assume throughout the paper that the covariance matrix $\Sigma = \mathbb{E}(X - \mu)(X - \mu)^T$ exists. Suppose that N independent, identically distributed samples X_1, \dots, X_N drawn from the distribution of X are available, and one wishes to estimate the mean vector μ . An estimator is simply a function of the data that we denote by $\hat{\mu}_N = \hat{\mu}_N(X_1, \dots, X_N)$.

There are many possible ways of measuring the quality of an estimator. The classical statistical literature tended to focus on risk measures such as the mean squared error $\mathbb{E}\|\hat{\mu}_N - \mu\|^2$. (Here, and in the rest of the paper, $\|\cdot\|$ denotes the Euclidean norm in \mathbb{R}^d , $S^{d-1} = \{v \in \mathbb{R}^d : \|v\| = 1\}$ denotes the Euclidean sphere in \mathbb{R}^d and $\langle \cdot, \cdot \rangle$ is the usual inner product in \mathbb{R}^d .) In this case, the sample mean $\bar{\mu}_N = (1/N) \sum_{i=1}^N X_i$ has a mean squared error equal to $\text{Tr}(\Sigma)/N$ [where $\text{Tr}(\Sigma)$ denotes the trace of the covariance matrix] and, even though this estimator is not necessarily optimal even for standard normal vectors—by “Stein’s paradox” (see [10])—the order of magnitude of the error cannot be improved in general.

The situation is quite different when one is interested in minimizing the value r that satisfies

$$\mathbb{P}\{\|\hat{\mu}_N - \mu\| > r\} \leq \delta$$

for some given $\delta > 0$. While one may always take $r = \sqrt{\text{Tr}(\Sigma)/(N\delta)}$ for the sample mean, much better dependence on δ may be achieved if the distribution is sufficiently light tailed. For example, if X has a multivariate normal distribution with

Received February 2017; revised July 2017.

¹Supported by the Spanish Ministry of Economy and Competitiveness Grant MTM2015-67304-P and FEDER, EU.

²Supported in part by the Israel Science Foundation.

MSC2010 subject classifications. Primary 62J02, 62G08; secondary 60G25.

Key words and phrases. Mean estimation, robust estimation, sub-Gaussian inequalities.

mean μ and covariance matrix Σ , then the sample mean $\bar{\mu}_N$ is also multivariate normal with mean μ and covariance matrix $(1/N)\Sigma$ and, therefore, for $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$(1.1) \quad \|\bar{\mu}_N - \mu\| \leq \sqrt{\frac{\text{Tr}(\Sigma)}{N}} + \sqrt{\frac{2\lambda_{\max} \log(1/\delta)}{N}},$$

where λ_{\max} denotes the largest eigenvalue of Σ (see Hanson and Wright [7]). Similar bounds may be proven for the performance of the sample mean if X has a sub-Gaussian distribution in the sense that for all unit vectors $v \in S^{d-1}$,

$$\mathbb{E} \exp(\lambda \langle v, X - \mathbb{E}X \rangle) \leq \exp(c\lambda^2 \langle v, \Sigma v \rangle)$$

for some constant c .

However, when the distribution is not necessarily sub-Gaussian and is possibly heavy tailed, one cannot expect such a sub-Gaussian behavior of the sample mean. Thus, when is it not reasonable to assume a sub-Gaussian distribution and heavy tails may be a concern, the sample mean is a risky choice. Indeed, alternative estimators have been constructed to achieve better performance.

The one-dimensional case (i.e., $d = 1$) is quite well understood; see Catoni [4] and Devroye, Lerasle, Lugosi and Oliveira [6] for recent accounts. The so-called *median-of-means* estimator is a simple and powerful univariate estimator with essentially optimal performance. This estimate was introduced independently in various papers; see Nemirovsky and Yudin [17], Jerrum, Valiant and Vazirani [11], Alon, Matias and Szegedy [1]. The median-of-means estimator partitions the data into $k < N$ blocks of size $m \approx N/k$ each, computes the sample mean within each block and outputs their median. One may easily show (see, e.g., Hsu [8]) that, for any $\delta \in (0, 1)$ if $k = \lceil 8 \log(1/\delta) \rceil$, then the resulting estimator $\hat{\mu}_N^{(\delta)}$ satisfies that, with probability at least $1 - \delta$,

$$(1.2) \quad |\hat{\mu}_N^{(\delta)} - \mu| \leq 8\sigma \sqrt{\frac{\log(2/\delta)}{N}},$$

where σ^2 denotes the variance of X . In other words, in the one-dimensional case, the median-of-means estimator achieves a sub-Gaussian performance under the only condition that the variance of X exists.

The median-of-means estimator has been extended to the multivariate case by replacing the median by its natural multivariate extension, the so-called “geometric (or spatial) median” (i.e., the point that minimizes the sum of the Euclidean distances to the sample means within each block); see Lerasle and Oliveira [14], Hsu and Sabato [9], Minsker [16]. In particular, Minsker proves that for each $\delta \in (0, 1)$ this generalization of the median-of-means estimator $\tilde{\mu}_N^{(\delta)}$ is such that, with probability at least $1 - \delta$,

$$(1.3) \quad \|\tilde{\mu}_N^{(\delta)} - \mu\| \leq C \sqrt{\frac{\text{Tr}(\Sigma) \log(1/\delta)}{N}},$$

where C is a universal constant. This bound holds under the only assumption that the covariance matrix exists. However, it does not quite achieve a sub-Gaussian performance bound that resembles (1.1).

Joly, Lugosi and Oliveira [12] made an attempt to construct a mean estimator with a sub-Gaussian behavior for a large class of distributions. They prove that there exists a mean estimator $\widehat{\mu}_n^{(\delta)}$ such that, if the distribution satisfies that for all $v \in S^{d-1}$

$$\mathbb{E}[(X - \mu, v)^4] \leq K (v, \Sigma v)^2$$

for some constant K , then for all $N \geq CK \log d(d + \log(1/\delta))$, with probability at least $1 - \delta$,

$$(1.4) \quad \|\widehat{\mu}_N^{(\delta)} - \mu\| \leq C \left(\sqrt{\frac{\text{Tr}(\Sigma)}{N}} + \sqrt{\frac{\lambda_{\max} \log(\delta^{-1} \log d)}{N}} \right),$$

where again C is a universal constant. This bound resembles the sub-Gaussian inequality (1.1). However, there are various caveats: the additional fourth-moment assumption, the requirement that $N = \Omega(d \log d)$, and, to a lesser extent, the extra $\log \log d$ term in the bound seems suboptimal.

The main result of this paper is that there exists a mean estimator that achieves purely sub-Gaussian performance under the minimal condition that the covariance matrix exists. More precisely, we prove the existence of a mean estimator $\widehat{\mu}_N^{(\delta)}$ such that, for all distributions with a finite second moment, for all N , with probability at least $1 - \delta$,

$$\|\widehat{\mu}_N^{(\delta)} - \mu\| \leq C \left(\sqrt{\frac{\text{Tr}(\Sigma)}{N}} + \sqrt{\frac{\lambda_{\max} \log(2/\delta)}{N}} \right)$$

for an explicit universal constant C .

The proposed estimator may be interpreted as a multivariate median-of-means estimate but with a new notion of a multivariate median which may be interesting in its own right. The construction of the new estimator is inspired by the technique of “median-of-means tournament,” put forward by the authors in [15].

In the next section, we present the proposed estimator and the performance bound. In Section 3, we present the proofs. We finish the paper by remarks about the computation of the estimator.

2. The estimator. Here, we introduce a mean estimator with a sub-Gaussian performance for all distributions whose covariance matrix exists. Recall that we are given an i.i.d. sample X_1, \dots, X_n of random vectors in \mathbb{R}^d . As in the case of the median-of-means estimator, we start by partitioning the set $\{1, \dots, n\}$ into k blocks B_1, \dots, B_k , each of size $|B_j| \geq m \stackrel{\text{def.}}{=} \lfloor n/k \rfloor$, where k is a parameter of the estimator whose value depends on the desired confidence level, as specified

below. In order to simplify the presentation, in the rest of the paper, without loss of generality, we assume that n is divisible by k and, therefore, $|B_j| = m$ for all $j = 1, \dots, k$.

Define the sample mean within each block by

$$Z_j = \frac{1}{m} \sum_{i \in B_j} X_i.$$

For each $a \in \mathbb{R}^d$, let

$$(2.1) \quad T_a = \{x \in \mathbb{R}^d : \exists J \subset [k] : |J| > k/2 \text{ such that} \\ \text{for all } j \in J, \|Z_j - x\| \leq \|Z_j - a\|\}$$

and define the mean estimator by

$$\hat{\mu}_n \in \underset{a \in \mathbb{R}^d}{\operatorname{argmin}} \operatorname{radius}(T_a),$$

where $\operatorname{radius}(T_a) = \sup_{x \in T_a} \|x - a\|$. Thus, $\hat{\mu}_n$ is chosen to minimize, over all $a \in \mathbb{R}^d$, the radius of the set T_a defined as the set of points $x \in \mathbb{R}^d$ for which $\|Z_j - x\| \leq \|Z_j - a\|$ for the majority of the blocks. If there are several minimizers, one may pick any one of them.

Note that the minimum is always achieved. This follows from the fact that $\operatorname{radius}(T_a)$ is a continuous function of a (since, for each a , T_a is the intersection of a finite union of closed balls, and the centers and radii of the closed balls are continuous in a).

One may interpret $\operatorname{argmin}_{a \in \mathbb{R}^d} \operatorname{radius}(T_a)$ as a new multivariate notion of the median of Z_1, \dots, Z_k . Indeed, when $d = 1$, it is a particular choice of the median and the proposed estimator coincides with the median-of-means estimator.

The main result of this paper is the following performance bound.

THEOREM 1. *Let $\delta \in (0, 1)$ and consider the mean estimator $\hat{\mu}_n$ with parameter $k = \lceil 200 \log(2/\delta) \rceil$. If X_1, \dots, X_n are i.i.d. random vectors in \mathbb{R}^d with mean $\mu \in \mathbb{R}^d$ and covariance matrix Σ , then for all n , with probability at least $1 - \delta$,*

$$\|\hat{\mu}_n - \mu\| \leq \max\left(960\sqrt{\frac{\operatorname{Tr}(\Sigma)}{n}}, 240\sqrt{\frac{\lambda_{\max} \log(2/\delta)}{n}}\right).$$

Thus, the proposed estimator achieves a purely sub-Gaussian performance under minimal conditions. Just like in the case of the median-of-means estimator for the univariate case, the estimator depends on the desired level of confidence δ . As it is shown in [6], such a dependence cannot be avoided without imposing additional conditions on the distribution. However, following the route laid down in [6], one may construct sub-Gaussian estimators that work for a wide range of confidence

levels simultaneously under more assumptions on the distribution. Since this issue is beyond the scope of this paper, it will not be pursued further here.

Just like Minsker’s bound (1.3)—but unlike the bound (1.4)—the performance bound of Theorem 1 is “infinite-dimensional” in the sense that the bound does not depend on the dimension d explicitly. Indeed, the same estimator may be defined for Hilbert-space valued random vectors and Theorem 1 remains valid as long as $\text{Tr}(\Sigma) = \mathbb{E}\|X - \mu\|^2$ is finite.

Theorem 1 is an outcome of the following observation which is of interest in its own right on the geometry of a typical collection $\{X_1, \dots, X_n\}$.

THEOREM 2. *Using the same notation as above and setting*

$$r = \max\left(960\sqrt{\frac{\text{Tr}(\Sigma)}{n}}, 240\sqrt{\frac{\lambda_{\max} \log(2/\delta)}{n}}\right),$$

with probability at least $1 - \delta$, for any $a \in \mathbb{R}^d$ such that $\|a - \mu\| \geq r$, one has $\|Z_j - a\| > \|Z_j - \mu\|$ for more than $k/2$ indices j .

Theorem 2 implies that for a “typical” collection X_1, \dots, X_n , μ is closer to a majority of the Z_j ’s when compared to any $a \in \mathbb{R}^d$ that is sufficiently far from μ . Obviously, for an arbitrary collection $x_1, \dots, x_n \subset \mathbb{R}^d$ such a point need not exist, and it is rather surprising that for a typical i.i.d. configuration, this property is satisfied by μ .

The fact that Theorem 2 implies Theorem 1 is straightforward. Indeed, the definition of $\hat{\mu}_n$ and Theorem 2 imply that, with probability at least $1 - \delta$, $\text{radius}(T_{\hat{\mu}_n}) \leq \text{radius}(T_\mu) \leq r$. Since either $\mu \in T_{\hat{\mu}_n}$ or $\hat{\mu} \in T_\mu$, we must have $\|\hat{\mu}_n - \mu\| \leq r$, as required.

We do not claim that the values of the constants appearing in Theorem 1 are optimal. They were obtained with the goal of making the proof transparent, nothing more, and it is likely that they may be improved by more careful calculations.

The proof of Theorem 2 is based on the idea of “median-of-means tournaments,” which was introduced by Lugosi and Mendelson [15], in the context of regression function estimation.

3. Proof. The proof of Theorem 2 is based on the following idea. The mean μ is the minimizer of the function $f(x) = \mathbb{E}\|X - \mu\|^2$. A possible approach is to use the available data to guess, for any pair $a, b \in \mathbb{R}^d$, whether $f(a) < f(b)$. To this end, we may set up a “tournament” as follows.

Recall that $[n]$ is partitioned into k disjoint blocks B_1, \dots, B_k of size $m = n/k$. For $a, b \in \mathbb{R}^d$, we say that a *defeats* b if

$$\frac{1}{m} \sum_{i \in B_j} (\|X_i - b\|^2 - \|X_i - a\|^2) > 0$$

on more than $k/2$ blocks B_j . The main technical lemma is the following.

LEMMA 1. Let $\delta \in (0, 1)$, $k = \lceil 200 \log(2/\delta) \rceil$, and define

$$r = \max\left(960\sqrt{\frac{\text{Tr}(\Sigma)}{n}}, 240\sqrt{\frac{\lambda_{\max} \log(2/\delta)}{n}}\right).$$

With probability at least $1 - \delta$, μ defeats all $b \in \mathbb{R}^d$ such that $\|b - \mu\| \geq r$.

PROOF. Note that

$$\begin{aligned} \|X_i - b\|^2 - \|X_i - \mu\|^2 &= \|X_i - \mu + \mu - b\|^2 - \|X_i - \mu\|^2 \\ &= -2\langle X_i - \mu, b - \mu \rangle + \|b - \mu\|^2, \end{aligned}$$

set $\bar{X} = X - \mu$ and put $v = b - \mu$. Thus, for a fixed b that satisfies $\|b - \mu\| \geq r$, μ defeats b if

$$-\frac{2}{m} \sum_{i \in B_j} \langle \bar{X}_i, v \rangle + \|v\|^2 > 0$$

on the majority of blocks B_j .

Therefore, to prove our claim we need that, with probability at least $1 - \delta$, for every $v \in \mathbb{R}^d$ with $\|v\| \geq r$,

$$(3.1) \quad -\frac{2}{m} \sum_{i \in B_j} \langle \bar{X}_i, v \rangle + \|v\|^2 > 0$$

for more than $k/2$ blocks B_j . Clearly, it suffices to show that (3.1) holds when $\|v\| = r$.

Consider a fixed $v \in \mathbb{R}^d$ with $\|v\| = r$. By Chebyshev’s inequality, with probability at least $9/10$,

$$\left| \frac{1}{m} \sum_{i \in B_j} \langle \bar{X}_i, v \rangle \right| \leq \sqrt{10} \sqrt{\frac{\mathbb{E} \langle \bar{X}, v \rangle^2}{m}} \leq \sqrt{10} \|v\| \sqrt{\frac{\lambda_{\max}}{m}},$$

where recall that λ_{\max} is the largest eigenvalue of the covariance matrix of X . Thus, if

$$(3.2) \quad r = \|v\| \geq 4\sqrt{10} \sqrt{\frac{\lambda_{\max}}{m}}$$

then with probability at least $9/10$,

$$(3.3) \quad -\frac{2}{m} \sum_{i \in B_j} \langle \bar{X}_i, v \rangle \geq \frac{-r^2}{2}.$$

Applying a standard binomial tail estimate, we see that (3.3) holds for a single v with probability at least $1 - \exp(-k/50)$ on at least $8/10$ of the blocks B_j .

Now we need to extend the above from a fixed vector v to all vectors with norm r . In order to show that (3.3) holds simultaneously for all $v \in r \cdot S^{d-1}$ on at least 7/10 of the blocks B_j , we first consider a maximal ε -separated set $V_1 \subset r \cdot S^{d-1}$ with respect to the $L_2(X)$ norm. In other words, V_1 is a subset of $r \cdot S^{d-1}$ of maximal cardinality such that for all $v_1, v_2 \in V_1$, $\|v_1 - v_2\|_{L_2(X)} = \langle v_1 - v_2, \Sigma(v_1 - v_2) \rangle^{1/2} \geq \varepsilon$. We may estimate this cardinality by the ‘‘dual Sudakov’’ inequality (Proposition 1 in the Appendix), which implies that the cardinality of V_1 is bounded by

$$\log(|V_1|/2) \leq \frac{1}{32} \left(\frac{\mathbb{E}[\langle G, \Sigma G \rangle^{1/2}]}{\varepsilon/r} \right)^2,$$

where G is a standard normal vector in \mathbb{R}^d . Notice that for any $a \in \mathbb{R}^d$, $\mathbb{E}_X \langle a, X \rangle^2 = \langle a, \Sigma a \rangle$ and, therefore,

$$\begin{aligned} \mathbb{E}[\langle G, \Sigma G \rangle^{1/2}] &= \mathbb{E}_G[(\mathbb{E}_X[\langle G, \bar{X} \rangle^2])^{1/2}] \leq (\mathbb{E}_X \mathbb{E}_G[\langle G, \bar{X} \rangle^2])^{1/2} \\ &= (\mathbb{E}[\|\bar{X}\|^2])^{1/2} = \sqrt{\text{Tr}(\Sigma)}. \end{aligned}$$

Hence, by setting

$$(3.4) \quad \varepsilon = 2r \left(\frac{1}{k} \text{Tr}(\Sigma) \right)^{1/2},$$

we have $|V_1| \leq 2e^{k/100}$ and thus, by the union bound, with probability at least $1 - 2e^{-k/100} \geq 1 - \delta/2$, (3.3) holds for all $v \in V_1$ on at least 8/10 of the blocks B_j .

Next, we check that property (3.1) holds simultaneously for all x with $\|x\| = r$ on at least 7/10 of the blocks B_j .

For every $x \in r \cdot S^{d-1}$, let v_x be the nearest element to x in V_1 with respect to the $L_2(X)$ norm. It suffices to show that, with probability at least $1 - \exp(-k/200) \geq 1 - \delta/2$,

$$(3.5) \quad \sup_{x \in r \cdot S^{d-1}} \frac{1}{k} \sum_{j=1}^k \mathbb{1}_{\{|m^{-1} \sum_{i \in B_j} \langle \bar{X}_i, x - v_x \rangle| \geq r^2/4\}} \leq \frac{1}{10}.$$

Indeed, on that event it follows that for every $x \in r \cdot S^{d-1}$, on at least 7/10 of the coordinate blocks B_j , both

$$-\frac{2}{m} \sum_{i \in B_j} \langle \bar{X}_i, v_x \rangle \geq \frac{-r^2}{2} \quad \text{and} \quad 2 \left| \frac{1}{m} \sum_{i \in B_j} \langle \bar{X}_i, x \rangle - \frac{1}{m} \sum_{i \in B_j} \langle \bar{X}_i, v_x \rangle \right| < \frac{r^2}{2}$$

hold, and hence, on those blocks, $-\frac{2}{m} \sum_{i \in B_j} \langle \bar{X}_i, x \rangle + r^2 > 0$ as required.

It remains to prove (3.5). Observe that

$$\frac{1}{k} \sum_{j=1}^k \mathbb{1}_{\{|m^{-1} \sum_{i \in B_j} \langle \bar{X}_i, x - v_x \rangle| \geq r^2/4\}} \leq \frac{4}{r^2} \frac{1}{k} \sum_{j=1}^k \left| \frac{1}{m} \sum_{i \in B_j} \langle \bar{X}_i, x - v_x \rangle \right|.$$

Since $\|x - v_x\|_{L_2(X)} = (\mathbb{E}\langle X, x - v_x \rangle^2)^{1/2} \leq \varepsilon$, it follows that for every j

$$\mathbb{E} \left| \frac{1}{m} \sum_{i \in B_j} \langle \bar{X}_i, x - v_x \rangle \right| \leq \sqrt{\frac{\mathbb{E}[\langle \bar{X}, x - v_x \rangle^2]}{m}} \leq \frac{\varepsilon}{\sqrt{m}}$$

and, therefore,

$$\begin{aligned} & \mathbb{E} \sup_{x \in r \cdot S^{d-1}} \frac{1}{k} \sum_{j=1}^k \mathbb{1}_{\{|m^{-1} \sum_{i \in B_j} \langle \bar{X}_i, x - v_x \rangle| \geq r^2/4\}} \\ & \leq \frac{4}{r^2} \mathbb{E} \sup_{x \in r \cdot S^{d-1}} \frac{1}{k} \sum_{j=1}^k \left(\left| \frac{1}{m} \sum_{i \in B_j} \langle \bar{X}_i, x - v_x \rangle \right| - \mathbb{E} \left| \frac{1}{m} \sum_{i \in B_j} \langle \bar{X}_i, x - v_x \rangle \right| \right) \\ & \quad + \frac{4\varepsilon}{r^2 \sqrt{m}} \\ & \stackrel{\text{def.}}{=} (A) + (B). \end{aligned}$$

To bound (B), note that, by (3.4),

$$\frac{4\varepsilon}{r^2 \sqrt{m}} = 8 \left(\frac{\text{Tr}(\Sigma)}{n} \right)^{1/2} \cdot \frac{1}{r} \leq \frac{1}{60}$$

provided that

$$r \geq 480 \left(\frac{\text{Tr}(\Sigma)}{n} \right)^{1/2}.$$

We may bound (A) by standard techniques of empirical processes such as symmetrization, contraction for Rademacher averages and de-symmetrization. Indeed, let $\sigma_1, \dots, \sigma_n$ be independent Rademacher random variables (i.e., $\mathbb{P}\{\sigma_i = 1\} = \mathbb{P}\{\sigma_i = -1\} = 1/2$), independent of the X_i . Then

$$\begin{aligned} (A) & \leq \frac{8}{r^2} \mathbb{E} \sup_{x \in r \cdot S^{d-1}} \frac{1}{k} \sum_{j=1}^k \sigma_j \left| \frac{1}{m} \sum_{i \in B_j} \langle \bar{X}_i, x - v_x \rangle \right| \\ & \quad \text{(by the first inequality of Proposition 2 below)} \\ & \leq \frac{8}{r^2} \mathbb{E} \sup_{x \in r \cdot S^{d-1}} \left| \frac{1}{k} \sum_{j=1}^k \sigma_j \frac{1}{m} \sum_{i \in B_j} \langle \bar{X}_i, x - v_x \rangle \right| \\ & \quad \text{(by Proposition 3 below)} \\ & \leq \frac{16}{r^2} \mathbb{E} \sup_{x \in r \cdot S^{d-1}} \left| \frac{1}{n} \sum_{i=1}^n \langle \bar{X}_i, x - v_x \rangle \right| \\ & \quad \text{(by the second inequality of Proposition 2)} \end{aligned}$$

$$\begin{aligned} &\leq \frac{32}{r} \mathbb{E} \sup_{\{t: \|t\| \leq 1\}} \left| \frac{1}{n} \sum_{i=1}^n \langle \bar{X}_i, t \rangle \right| \\ &\quad \text{(noting that } \|x - v_x\| \leq 2r) \\ &\leq \frac{32}{r} \cdot \frac{\mathbb{E} \|\bar{X}\|}{\sqrt{n}} = \frac{32}{r} \left(\frac{\text{Tr}(\Sigma)}{n} \right)^{1/2} \leq \frac{1}{30} \end{aligned}$$

provided that $r \geq 960 \left(\frac{\text{Tr}(\Sigma)}{n} \right)^{1/2}$.

Thus, for

$$Y = \sup_{x \in r \cdot S^{d-1}} \frac{1}{k} \sum_{j=1}^k \mathbb{1}_{\{|m^{-1} \sum_{i \in B_j} \langle \bar{X}_i, x - v_x \rangle| \geq r^2/4\}},$$

we have proved that $\mathbb{E}Y \leq 1/60 + 1/30 = 1/20$. Finally, in order to prove (3.5), it suffices to prove that, $\mathbb{P}\{Y > \mathbb{E}Y + 1/20\} \leq e^{-k/200}$, which follows from the bounded differences inequality (see, e.g., [3], Theorem 6.2). \square

PROOF OF THEOREM 2. Theorem 2 is easily derived from Lemma 1. Fix a block B_j , and recall that $Z_j = \frac{1}{m} \sum_{i \in B_j} X_i$. Let $a, b \in \mathbb{R}^d$. Then

$$\begin{aligned} \frac{1}{m} \sum_{i \in B_j} (\|X_i - a\|^2 - \|X_i - b\|^2) &= \frac{1}{m} \sum_{i \in B_j} (\|X_i - b - (a - b)\|^2 - \|X_i - b\|^2) \\ &= -\frac{2}{m} \sum_{i \in B_j} \langle X_i - b, a - b \rangle + \|a - b\|^2 = (*). \end{aligned}$$

Observe that $-\frac{2}{m} \sum_{i \in B_j} \langle X_i - b, a - b \rangle = -2 \langle \frac{1}{m} \sum_{i \in B_j} X_i - b, a - b \rangle = -2 \langle Z_j - b, a - b \rangle$, and thus

$$\begin{aligned} (*) &= -2 \langle Z_j - b, a - b \rangle + \|a - b\|^2 \\ &= -2 \langle Z_j - b, a - b \rangle + \|a - b\|^2 + \|Z_j - b\|^2 - \|Z_j - b\|^2 \\ &= \|Z_j - b - (a - b)\|^2 - \|Z_j - b\|^2 = \|Z_j - a\|^2 - \|Z_j - b\|^2. \end{aligned}$$

Therefore, $(*) > 0$ (i.e., b defeats a on block B_j) if and only if $\|Z_j - a\| > \|Z_j - b\|$.

Recall that Lemma 1 states that, with probability at least $1 - \delta$, if $\|a - \mu\| \geq r$ then on more than $k/2$ blocks B_j , $\frac{1}{m} \sum_{i \in B_j} (\|X_i - a\|^2 - \|X_i - \mu\|^2) > 0$, which, by the above argument, is the same as saying that for at least $k/2$ indices j , $\|Z_j - a\| > \|Z_j - \mu\|$. \square

4. Computational considerations. The problem of computing various notions of multivariate medians has been thoroughly studied in computational geometry and we refer to Aloupis [2] for a survey on this topic. For example, computing the geometric median and, therefore, the multivariate median-of-means estimator proposed by Hsu and Sabato [9] and Minsker involves solving a convex optimization problem. Thus, the geometric median may be approximated efficiently; see [5] for the most recent result and for the rich history of the problem.

In contrast, efficiently computing, or even approximating, the multivariate median proposed in this paper appears to be a nontrivial challenge.

A possible approach for computing a mean estimator that approximates $\hat{\mu}_n$ is based on a variant of a coordinate descent algorithm that works roughly as follows: starting with an arbitrary line in \mathbb{R}^d , one may discretize, with mesh $O(r)$, the segment on the line that supports the convex hull of Z_1, \dots, Z_k . Then one uses pairwise comparisons of the discretized values, using the median-of-means estimate, to find a point that defeats every other candidate on the line that is at least distance $2r$ apart from it. (With a minor adjustment of our arguments above, one may prove that such a point always exists.) Then take a line that is orthogonal to the first line and contains the “winner” and repeat the search on that line. Continue for d steps. One may prove that the point $\tilde{\mu}_n$ obtained at the final step is such that, with probability at least $1 - \delta$, $\|\tilde{\mu}_n - \mu\|_\infty \leq Cr$ for a numerical constant C . This algorithm runs in time quadratic in $1/r$ and linear in d but unfortunately it only guarantees closeness to the true mean in the ℓ_∞ sense. If one replaces orthogonal lines by random ones and keeps repeating the procedure, one eventually achieves the desired guarantee in the Euclidean distance. However, one needs to consider exponentially many (in d) directions to approach μ with the desired precision. Note that such algorithms use r as an input parameter. Naturally, the value of r is not known but the algorithm is guaranteed to work well as long as the true value of r is larger than the prior guess.

Another possibility is to start with computing the geometric median $\tilde{\mu}^{(\delta)}$ of the Z_j . By (1.3), one may now restrict search to a ball of radius at most $r\sqrt{\log(1/\delta)}$. By exhaustively searching through this ball (after appropriately discretizing), one finds an estimate with the desired properties in additional time of order $\log^d(1/\delta)$. However, this is surely unrealistic in most interesting cases.

We leave the question of efficiently computing the proposed mean estimate (or another one with sub-Gaussian performance guarantees) as an interesting research problem.

APPENDIX: SOME TECHNICAL TOOLS

Here, we list some of the standard tools of geometric analysis and empirical process theory used in the proofs.

We start with the so-called “dual-Sudakov” inequality; see [13] and also [19] for a version with the specified constant below. For a convex body $K \subset \mathbb{R}^d$ (i.e.,

a centrally-symmetric convex set with a nonempty interior), we denote by $N(K)$ the smallest number of translates of K needed to cover the Euclidean unit ball $\{x \in \mathbb{R}^d : \|x\| \leq 1\}$.

PROPOSITION 1 (Dual Sudakov inequality). *Let K be a convex body in \mathbb{R}^d . Then*

$$\sqrt{\log(N(K)/2)} \leq \frac{1}{\sqrt{32}} \mathbb{E} \|G\|_K,$$

where G is standard Gaussian vector in \mathbb{R}^d and $\|\cdot\|_K$ denotes the norm whose unit ball is K .

We also need the following symmetrization inequalities; see, for example, [18], Lemma 2.3.6.

PROPOSITION 2 (Symmetrization inequalities). *Let X_1, \dots, X_n be i.i.d. random vectors taking values in \mathbb{R}^d . Let \mathcal{F} be a class of real-valued functions defined on \mathbb{R}^d . Let $\sigma_1, \dots, \sigma_n$ be independent Rademacher random variables (i.e., $\mathbb{P}\{\sigma_i = 1\} = \mathbb{P}\{\sigma_i = -1\} = 1/2$), independent of the X_i . Then*

$$\mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (f(X_i) - \mathbb{E} f(X_i)) \leq 2 \mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(X_i).$$

Moreover, if $\mathbb{E} f(X_i) = 0$ for all $f \in \mathcal{F}$, then

$$\mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(X_i) \leq 2 \mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(X_i).$$

The following contraction lemma for Rademacher averages may be found in [13].

PROPOSITION 3 (Contraction lemma). *Let X_1, \dots, X_n be i.i.d. random vectors taking values in \mathbb{R}^d . Let \mathcal{F} be a class of real-valued functions defined on \mathbb{R}^d . Let $\sigma_1, \dots, \sigma_n$ be independent Rademacher random variables, independent of the X_i . If $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is a function with $\phi(0) = 0$ and Lipschitz constant L , then*

$$\mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \phi(f(X_i)) \leq L \cdot \mathbb{E} \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(X_i).$$

REFERENCES

- [1] ALON, N., MATIAS, Y. and SZEGEDY, M. (1999). The space complexity of approximating the frequency moments. *J. Comput. System Sci.* **58** 137–147. [MR1688610](#)
- [2] ALOUPIS, G. (2006). Geometric measures of data depth. In *Data Depth: Robust Multivariate Analysis, Computational Geometry and Applications*. DIMACS Ser. Discrete Math. Theoret. Comput. Sci. **72** 147–158. Amer. Math. Soc., Providence, RI. [MR2343118](#)

- [3] BOUCHERON, S., LUGOSI, G. and MASSART, P. (2013). *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford Univ. Press, Oxford.
- [4] CATONI, O. (2012). Challenging the empirical mean and empirical variance: A deviation study. *Ann. Inst. Henri Poincaré Probab. Stat.* **48** 1148–1185.
- [5] COHEN, M. B., LEE, Y. T., MILLER, G., PACHOCKI, J. and SIDFORD, A. (2016). Geometric median in nearly linear time. In *STOC'16—Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing* 9–21. ACM, New York. [MR3536551](#)
- [6] DEVROYE, L., LERASLE, M., LUGOSI, G. and OLIVEIRA, R. I. (2016). Sub-Gaussian mean estimators. *Ann. Statist.* **44** 2695–2725. [MR3576558](#)
- [7] HANSON, D. L. and WRIGHT, F. T. (1971). A bound on tail probabilities for quadratic forms in independent random variables. *Ann. Math. Stat.* **42** 1079–1083. [MR0279864](#)
- [8] HSU, D. (2010). Robust statistics. Available at <http://www.inherentuncertainty.org/2010/12/robust-statistics.html>.
- [9] HSU, D. and SABATO, S. (2016). Loss minimization and parameter estimation with heavy tails. *J. Mach. Learn. Res.* **17** Paper No. 18. [MR3491112](#)
- [10] JAMES, W. and STEIN, C. (1961). Estimation with quadratic loss. In *Proc. 4th Berkeley Sympos. Math. Statist. and Prob., Vol. I* 361–379. Univ. California Press, Berkeley, CA. [MR0133191](#)
- [11] JERRUM, M. R., VALIANT, L. G. and VAZIRANI, V. V. (1986). Random generation of combinatorial structures from a uniform distribution. *Theoret. Comput. Sci.* **43** 169–188. [MR0855970](#)
- [12] JOLY, E., LUGOSI, G. and OLIVEIRA, R. I. (2016). On the estimation of the mean of a random vector. Preprint.
- [13] LEDOUX, M. and TALAGRAND, M. (1991). *Probability in Banach Space*. Springer, New York.
- [14] LERASLE, M. and OLIVEIRA, R. I. (2012). Robust empirical mean estimators. Available at [arXiv:1112.3914](https://arxiv.org/abs/1112.3914).
- [15] LUGOSI, G. and MENDELSON, S. (2016). Risk minimization by median-of-means tournaments. Preprint.
- [16] MINSKER, S. (2015). Geometric median and robust estimation in Banach spaces. *Bernoulli* **21** 2308–2335.
- [17] NEMIROVSKY, A. S. and YUDIN, D. B. (1983). Problem complexity and method efficiency in optimization.
- [18] VAN DER WAART, A. W. and WELLNER, J. A. (1996). *Weak Convergence and Empirical Processes*. Springer, New York.
- [19] VERSHYNIN, R. (2009). Lectures in geometric functional analysis. Available at <https://www.math.uci.edu/~rvershyn/papers/GFA-book.pdf>.

ICREA
 PG. LLUÍS COMPANYS 23
 08010 BARCELONA
 SPAIN
 AND
 DEPARTMENT OF ECONOMICS AND BUSINESS
 POMPEU FABRA UNIVERSITY
 BARCELONA
 SPAIN
 AND
 BARCELONA GSE
 RAMON TRIAS FARGAS, 25-27
 08005 BARCELONA
 SPAIN
 E-MAIL: gabor.lugosi@upf.edu

DEPARTMENT OF MATHEMATICS
 TECHNION, I.I.T.
 HAIFA
 ISRAEL
 AND
 MATHEMATICAL SCIENCES INSTITUTE
 THE AUSTRALIAN NATIONAL UNIVERSITY
 CANBERRA 0200
 AUSTRALIA
 E-MAIL: shahar@tx.technion.ac.il