

Mean estimation: median-of-means tournaments

Gábor Lugosi

ICREA, Pompeu Fabra University, BGSE

based on joint work with

Luc Devroye (McGill, Montreal)

Matthieu Lerasle (CNRS, Nice)

Roberto Imbuzeiro Oliveira (IMPA, Rio)

Shahar Mendelson (Technion and ANU)

estimating the mean

Given $\mathbf{X}_1, \dots, \mathbf{X}_n$, a real i.i.d. sequence, estimate $\mu = \mathbb{E}\mathbf{X}_1$.

estimating the mean

Given $\mathbf{X}_1, \dots, \mathbf{X}_n$, a real i.i.d. sequence, estimate $\mu = \mathbb{E}\mathbf{X}_1$.

“Obvious” choice: empirical mean

$$\bar{\mu}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$$

estimating the mean

Given $\mathbf{X}_1, \dots, \mathbf{X}_n$, a real i.i.d. sequence, estimate $\mu = \mathbb{E}\mathbf{X}_1$.

“Obvious” choice: empirical mean

$$\bar{\mu}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$$

By the central limit theorem, if \mathbf{X} has a finite variance σ^2 ,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ \sqrt{n} |\bar{\mu}_n - \mu| > \sigma \sqrt{2 \log(2/\delta)} \right\} \leq \delta .$$

We would like **non-asymptotic inequalities** of a similar form.

estimating the mean

Given $\mathbf{X}_1, \dots, \mathbf{X}_n$, a real i.i.d. sequence, estimate $\mu = \mathbb{E}\mathbf{X}_1$.

“Obvious” choice: empirical mean

$$\bar{\mu}_n = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i$$

By the central limit theorem, if \mathbf{X} has a finite variance σ^2 ,

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ \sqrt{n} |\bar{\mu}_n - \mu| > \sigma \sqrt{2 \log(2/\delta)} \right\} \leq \delta .$$

We would like **non-asymptotic inequalities** of a similar form.

If the distribution is sub-Gaussian,

$\mathbb{E} \exp(\lambda(\mathbf{X} - \mu)) \leq \exp(\sigma^2 \lambda^2 / 2)$, then with probability at least $1 - \delta$,

$$|\bar{\mu}_n - \mu| \leq \sigma \sqrt{\frac{2 \log(2/\delta)}{n}} .$$

empirical mean–heavy tails

The empirical mean is computationally attractive.

Requires no a priori knowledge and automatically scales with σ .

If the distribution is not sub-Gaussian, we still have Chebyshev's inequality: w.p. $\geq 1 - \delta$,

$$|\bar{\mu}_n - \mu| \leq \sigma \sqrt{\frac{1}{n\delta}} .$$

Exponentially weaker bound. Especially hurts when many means are estimated simultaneously.

This is the best one can say. Catoni (2012) shows that for each δ there exists a distribution with variance σ such that

$$\mathbb{P} \left\{ |\bar{\mu}_n - \mu| \geq \sigma \sqrt{\frac{c}{n\delta}} \right\} \geq \delta .$$

median of means

A simple estimator is **median-of-means**. Goes back to Nemirovsky, Yudin (1983), Jerrum, Valiant, and Vazirani (1986), Alon, Matias, and Szegedy (2002).

$$\hat{\mu}_{MM} \stackrel{\text{def}}{=} \text{median} \left(\frac{1}{m} \sum_{t=1}^m \mathbf{x}_t, \dots, \frac{1}{m} \sum_{t=(k-1)m+1}^{km} \mathbf{x}_t \right)$$

median of means

A simple estimator is **median-of-means**. Goes back to Nemirovsky, Yudin (1983), Jerrum, Valiant, and Vazirani (1986), Alon, Matias, and Szegedy (2002).

$$\hat{\mu}_{MM} \stackrel{\text{def}}{=} \text{median} \left(\frac{1}{m} \sum_{t=1}^m \mathbf{x}_t, \dots, \frac{1}{m} \sum_{t=(k-1)m+1}^{km} \mathbf{x}_t \right)$$

Lemma

Let $\delta \in (0, 1)$, $k = 8 \log \delta^{-1}$ and $m = \frac{n}{8 \log \delta^{-1}}$. Then with probability at least $1 - \delta$,

$$|\hat{\mu}_{MM} - \mu| \leq \sigma \sqrt{\frac{32 \log(1/\delta)}{n}}$$

proof

By Chebyshev, each mean is within distance $\sigma\sqrt{4/m}$ of μ with probability $3/4$.

The probability that the median is not within distance $\sigma\sqrt{4/m}$ of μ is at most $\mathbb{P}\{\text{Bin}(k, 1/4) > k/2\}$ which is exponentially small in k .

median of means

- Sub-Gaussian deviations.
- Scales automatically with σ .
- Parameters depend on required confidence level δ .
- See Lerasle and Oliveira (2012), Hsu and Sabato (2013), Minsker (2014) for generalizations.
- Also works when the variance is infinite. If $\mathbb{E}[|\mathbf{X} - \mathbb{E}\mathbf{X}|^{1+\alpha}] = M$ for some $\alpha \leq 1$, then, with probability at least $1 - \delta$,

$$|\hat{\mu}_{MM} - \mu| \leq \left(8 \frac{(12M)^{1/\alpha} \ln(1/\delta)}{n} \right)^{\alpha/(1+\alpha)}$$

why sub-Gaussian?

Sub-Gaussian bounds are the best one can hope for when the variance is finite.

In fact, for any $M > 0, \alpha \in (0, 1], \delta > 2e^{-n/4}$, and mean estimator $\hat{\mu}_n$, there exists a distribution $\mathbb{E} [|X - \mathbb{E}X|^{1+\alpha}] = M$ such that

$$|\hat{\mu}_n - \mu| \geq \left(\frac{M^{1/\alpha} \ln(1/\delta)}{n} \right)^{\alpha/(1+\alpha)}.$$

Proof: The distributions $P_+(0) = 1 - p, P_+(c) = p$ and $P_-(0) = 1 - p, P_-(-c) = p$ are indistinguishable if all n samples are equal to 0.

why sub-Gaussian?

This shows **optimality of the median-of-means estimator** for all α .

It also shows that finite variance is necessary even for rate $n^{-1/2}$.

One cannot hope to get anything better than sub-Gaussian tails.

Catoni proved that **sample mean is optimal for the class of Gaussian distributions.**

multiple- δ estimators

Do there exist estimators that are sub-Gaussian simultaneously for all confidence levels?

An estimator is multiple- δ -sub-Gaussian for a class of distributions \mathcal{P} and δ_{\min} if for all $\delta \in [\delta_{\min}, 1)$, and all distributions in \mathcal{P} ,

$$|\hat{\mu}_n - \mu| \leq L\sigma \sqrt{\frac{\log(2/\delta)}{n}}.$$

multiple- δ estimators

Do there exist estimators that are sub-Gaussian simultaneously for all confidence levels?

An estimator is multiple- δ -sub-Gaussian for a class of distributions \mathcal{P} and δ_{\min} if for all $\delta \in [\delta_{\min}, 1)$, and all distributions in \mathcal{P} ,

$$|\hat{\mu}_n - \mu| \leq L\sigma \sqrt{\frac{\log(2/\delta)}{n}}.$$

The picture is more complex than before.

known variance

Given $0 < \sigma_1 \leq \sigma_2 < \infty$, define the class

$$\mathcal{P}_2^{[\sigma_1^2, \sigma_2^2]} = \{P : \sigma_1^2 \leq \sigma_P^2 \leq \sigma_2^2\}.$$

Let $R = \sigma_2/\sigma_1$.

known variance

Given $0 < \sigma_1 \leq \sigma_2 < \infty$, define the class

$$\mathcal{P}_2^{[\sigma_1^2, \sigma_2^2]} = \{P : \sigma_1^2 \leq \sigma_P^2 \leq \sigma_2^2\}.$$

Let $R = \sigma_2/\sigma_1$.

- If R is **bounded** then there exists a multiple- δ -sub-Gaussian estimator with $\delta_{\min} = 4e^{1-n/2}$;
- If R is **unbounded** then there is no multiple- δ -sub-Gaussian estimate for any L and $\delta_{\min} \rightarrow 0$.

A sharp distinction.

The exponentially small value of δ_{\min} is best possible.

construction of multiple- δ estimator

Reminiscent to Lepski's method of adaptive estimation.

For $k = 1, \dots, K = \log_2(1/\delta_{\min})$, use the median-of-means estimator to construct confidence intervals I_k such that

$$\mathbb{P}\{\mu \notin I_k\} \leq 2^{-k}.$$

(This is where knowledge of σ_2 and boundedness of R is used.)

Define

$$\hat{k} = \min \left\{ k : \bigcap_{j=k}^K I_j \neq \emptyset \right\}.$$

Finally, let

$$\hat{\mu}_n = \text{mid point of } \bigcap_{j=\hat{k}}^K I_j$$

proof

For any $k = 1, \dots, K$,

$$\mathbb{P}\{|\hat{\mu}_n - \mu| > |I_k|\} \leq \mathbb{P}\{\exists j \geq k : \mu \notin I_j\}$$

because if $\mu \in \cap_{j=k}^K I_j$, then $\cap_{j=k}^K I_j$ is non-empty and therefore $\hat{\mu}_n \in \cap_{j=k}^K I_j$.

But

$$\mathbb{P}\{\exists j \geq k : \mu \notin I_j\} \leq \sum_{j=k}^K \mathbb{P}\{\mu \notin I_j\} \leq 2^{1-k}$$

higher moments

For $\eta \geq 1$ and $\alpha \in (2, 3]$, define

$$\mathcal{P}_{\alpha, \eta} = \{P : \mathbb{E}|X - \mu|^\alpha \leq (\eta \sigma)^\alpha\} .$$

Then for some $C = C(\alpha, \eta)$ there exists a multiple- δ estimator with a constant L and $\delta_{\min} = e^{-n/C}$ for all sufficiently large n .

k -regular distributions

This follows from a more general result:

Define

$$p_-(j) = \mathbb{P} \left\{ \sum_{i=1}^j \mathbf{x}_i \leq j\mu \right\} \quad \text{and} \quad p_+(j) = \mathbb{P} \left\{ \sum_{i=1}^j \mathbf{x}_i \geq j\mu \right\} .$$

A distribution is k -regular if

$$\forall j \geq k, \min(p_+(j), p_-(j)) \geq 1/3.$$

For this class there exists a multiple- δ estimator with a constant L and $\delta_{\min} = e^{-n/k}$ for all n .

multivariate distributions

Let \mathbf{X} be a random vector taking values in \mathbb{R}^d with mean $\boldsymbol{\mu} = \mathbb{E}\mathbf{X}$ and covariance matrix $\boldsymbol{\Sigma} = \mathbb{E}(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T$.

Given an i.i.d. sample $\mathbf{X}_1, \dots, \mathbf{X}_n$, we want to estimate $\boldsymbol{\mu}$ that has **sub-Gaussian** performance.

multivariate distributions

Let \mathbf{X} be a random vector taking values in \mathbb{R}^d with mean $\boldsymbol{\mu} = \mathbb{E}\mathbf{X}$ and covariance matrix $\boldsymbol{\Sigma} = \mathbb{E}(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T$.

Given an i.i.d. sample $\mathbf{X}_1, \dots, \mathbf{X}_n$, we want to estimate $\boldsymbol{\mu}$ that has **sub-Gaussian** performance.

What is sub-Gaussian?

If \mathbf{X} has a multivariate Gaussian distribution, the sample mean $\bar{\boldsymbol{\mu}}_n = (\mathbf{1}/n) \sum_{i=1}^n \mathbf{X}_i$ satisfies, with probability at least $\mathbf{1} - \delta$,

$$\|\bar{\boldsymbol{\mu}}_n - \boldsymbol{\mu}\| \leq \sqrt{\frac{\text{Tr}(\boldsymbol{\Sigma})}{n}} + \sqrt{\frac{2\lambda_{\max} \log(1/\delta)}{n}},$$

Can one construct mean estimators with similar performance for a large class of distributions?

coordinate-wise median of means

Coordinate-wise median of means yields the bound:

$$\|\hat{\mu}_{MM} - \mu\| \leq K \sqrt{\frac{\text{Tr}(\Sigma) \log(d/\delta)}{n}}.$$

We can do better.

multivariate median of means

Hsu and Sabato (2013), Minsker (2015) extended the median-of-means estimate.

Minsker proposes an analogous estimate that uses the multivariate median

$$\text{Med}(\mathbf{x}_1, \dots, \mathbf{x}_N) = \underset{\mathbf{y} \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{i=1}^N \|\mathbf{y} - \mathbf{x}_i\| .$$

For this estimate, with probability at least $1 - \delta$,

$$\|\hat{\mu}_{MM} - \mu\| \leq K \sqrt{\frac{\operatorname{Tr}(\Sigma) \log(1/\delta)}{n}} .$$

No further assumption or knowledge of the distribution is required.

Computationally feasible.

Almost sub-Gaussian but not quite.

Dimension free.

median-of-means tournament

We propose a new estimator with a purely sub-Gaussian performance, without further conditions.

The mean μ is the minimizer of $f(x) = \mathbb{E}\|X - \mu\|^2$.

For any pair $a, b \in \mathbb{R}^d$, we try to guess whether $f(a) < f(b)$ and set up a “tournament”.

Partition the data points into k blocks of size $m = n/k$.

We say that a defeats b if

$$\frac{1}{m} \sum_{i \in B_j} \|X_i - a\|^2 < \frac{1}{m} \sum_{i \in B_j} \|X_i - b\|^2$$

on more than $k/2$ blocks B_j .

median-of-means tournament

Within each block compute

$$Y_j = \frac{1}{m} \sum_{i \in B_j} x_i .$$

Then a defeats b if

$$\|Y_j - a\| < \|Y_j - b\|$$

on more than $k/2$ blocks B_j .

Lemma. Let $k = \lceil 200 \log(2/\delta) \rceil$. With probability at least $1 - \delta$, μ defeats all $b \in \mathbb{R}^d$ such that $\|b - \mu\| \geq r$, where

$$r = \max \left(800 \left(\sqrt{\frac{\text{Tr}(\Sigma)}{n}}, 240 \sqrt{\frac{\lambda_{\max} \log(2/\delta)}{n}} \right) \right) .$$

sub-gaussian estimate

For each $\mathbf{a} \in \mathbb{R}^d$, define the set

$$S_{\mathbf{a}} = \left\{ \mathbf{x} \in \mathbb{R}^d : \text{such that } \mathbf{x} \text{ defeats } \mathbf{a} \right\}$$

Now define the mean estimator as

$$\hat{\mu}_N \in \underset{\mathbf{a} \in \mathbb{R}^d}{\operatorname{argmin}} \operatorname{radius}(S_{\mathbf{a}}) .$$

By the lemma, w.p. $\geq 1 - \delta$,

$$\operatorname{radius}(S_{\hat{\mu}_N}) \leq \operatorname{radius}(S_{\mu}) \leq r$$

and therefore

$$\|\hat{\mu}_n - \mu\| \leq r .$$

sub-gaussian performance

Theorem. Let $k = \lceil 200 \log(2/\delta) \rceil$. Then, with probability at least $1 - \delta$,

$$\|\hat{\mu}_n - \mu\| \leq r$$

where

$$r = \max \left(800 \left(\sqrt{\frac{\text{Tr}(\Sigma)}{n}}, 240 \sqrt{\frac{\lambda_{\max} \log(2/\delta)}{n}} \right) \right) .$$

- No other condition other than existence of Σ .
- “Infinite-dimensional” inequality: the same holds in Hilbert spaces.
- The constants are explicit but sub-optimal.

proof of lemma: sketch

Let $\bar{\mathbf{X}} = \mathbf{X} - \mu$ and $\mathbf{v} = \mathbf{b} - \mu$. Then μ defeats \mathbf{b} if

$$-\frac{1}{m} \sum_{i \in B_j} \langle \bar{\mathbf{X}}_i, \mathbf{v} \rangle + \|\mathbf{v}\|^2 > 0$$

on the majority of blocks B_j . We need to prove that this holds for all \mathbf{v} with $\|\mathbf{v}\| = r$.

Step 1: For a fixed \mathbf{v} , by Chebyshev, with probability at least **9/10**,

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

So by a binomial tail estimate, with probability at least $1 - \exp(-k/50)$, this holds on at least **8/10** of the blocks B_j .

proof sketch

Step 2: Now we take a minimal ϵ cover the set $\mathbf{r} \cdot \mathbf{S}^{d-1}$ with respect to the norm $\langle \mathbf{v}, \Sigma \mathbf{v} \rangle^{1/2}$.

This set has $< e^{k/100}$ points if

$$\epsilon = 5r \left(\frac{1}{k} \text{Tr}(\Sigma) \right)^{1/2},$$

so we can use the union bound over this ϵ -net.

Step 3: To extend to all points in $\mathbf{r} \cdot \mathbf{S}^{d-1}$, we need that, with probability at least $1 - \exp(-k/200)$,

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

This may be proved by standard techniques of empirical processes.

algorithmic challenge

Computing the proposed estimator is an interesting open problem.

Coordinate descent does not quite do the job—it only guarantees

$$\|\hat{\mu}_n - \mu\|_\infty \leq r.$$

regression function estimation

Consider the standard statistical supervised learning problem under the squared loss.

Let (\mathbf{X}, \mathbf{Y}) take values in $\mathcal{X} \times \mathbb{R}$.

The goal is to predict \mathbf{Y} , upon observing \mathbf{X} , by $f(\mathbf{X})$ for some $f : \mathcal{X} \rightarrow \mathbb{R}$.

We measure the quality of f by the risk

$$\mathbb{E}(f(\mathbf{X}) - \mathbf{Y})^2 .$$

We have access to a sample $\mathcal{D}_n = ((\mathbf{X}_1, \mathbf{Y}_1), \dots, (\mathbf{X}_n, \mathbf{Y}_n))$.

We choose \hat{f}_n from a fixed class of functions \mathcal{F} . The best function is

$$f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \mathbb{E}(f(\mathbf{X}) - \mathbf{Y})^2 .$$

regression function estimation

We measure performance by either the **mean squared error**

$$\|\hat{f}_n - f^*\|_{L_2}^2 = \mathbb{E}((\hat{f}_n(X) - f^*(X))^2 | \mathcal{D}_n)$$

or by the **excess risk**

$$R(\hat{f}_n) = \mathbb{E}((\hat{f}_n(X) - Y)^2 | \mathcal{D}_n) - \mathbb{E}(f^*(X) - Y)^2 .$$

regression function estimation

We measure performance by either the **mean squared error**

$$\|\hat{f}_n - f^*\|_{L_2}^2 = \mathbb{E}((\hat{f}_n(X) - f^*(X))^2 | \mathcal{D}_n)$$

or by the **excess risk**

$$R(\hat{f}_n) = \mathbb{E}((\hat{f}_n(X) - Y)^2 | \mathcal{D}_n) - \mathbb{E}(f^*(X) - Y)^2 .$$

A procedure achieves **accuracy** r with **confidence** $1 - \delta$ if

$$\mathbb{P} \left(\|\hat{f}_n - f^*\|_{L_2} \leq r \right) \geq 1 - \delta .$$

High accuracy and high confidence are conflicting requirements.

The **accuracy edge** is the smallest achievable accuracy with confidence $1 - \delta = 3/4$.

A quest with a long history has been to understand the tradeoff.

empirical risk minimization

The standard learning procedure is **empirical risk minimization (ERM)**:

$$\hat{f}_n = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{i=1}^n (f(X_i) - Y_i)^2 .$$

ERM achieves near-optimal accuracy/confidence tradeoff for well-behaved distributions.

The performance of ERM is now well understood.

It works well if both \mathbf{Y} and $\mathbf{f}(\mathbf{X})$ have sub-Gaussian tails (for all $\mathbf{f} \in \mathcal{F}$).

four complexity parameters

The performance of ERM depends on the intricate interplay between the geometry of \mathcal{F} and the distribution of (X, Y) . We assume that \mathcal{F} is **convex**.

Let $\mathcal{F}_{h,r} = \{f - h : f \in \mathcal{F}, \|f - h\|_{L_2} \leq r\}$ and let $\mathcal{M}(\mathcal{F}_{h,r}, \epsilon)$ be the ϵ -packing numbers.

For $\kappa, \eta > 0$, set

$$\lambda_{\mathbb{Q}}(\kappa, \eta) = \sup_{h \in \mathcal{F}} \inf \{r : \log \mathcal{M}(\mathcal{F}_{h,r}, \eta r) \leq \kappa^2 n\} .$$

Similarly, let

$$\lambda_{\mathbb{M}}(\kappa, \eta) = \sup_{h \in \mathcal{F}} \inf \{r : \log \mathcal{M}(\mathcal{F}_{h,r}, \eta r) \leq \kappa^2 n r^2\}$$

four complexity parameters

$$r_E(\kappa) = \sup_{h \in \mathcal{F}} \inf \left\{ r : \mathbb{E} \sup_{u \in \mathcal{F}_{h,r}} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i u(X_i) \right| \leq \kappa \sqrt{nr} \right\},$$

Finally, let

$$\begin{aligned} \bar{r}_{\mathbb{M}}(\kappa, h) \\ = \inf \left\{ r : \mathbb{E} \sup_{u \in \mathcal{F}_{h,r}} \left| \frac{1}{\sqrt{n}} \sum_{i=1}^n \epsilon_i u(X_i) \cdot (h(X_i) - Y_i) \right| \leq \kappa \sqrt{nr^2} \right\}. \end{aligned}$$

and

$$\tilde{r}_{\mathbb{M}}(\kappa, \sigma) = \sup_{h \in \mathcal{F}_Y^{(\sigma)}} \bar{r}_{\mathbb{M}}(\kappa, h)$$

where $\mathcal{F}_Y^{(\sigma)} = \{f \in \mathcal{F} : \|f(X) - Y\|_{L_2} \leq \sigma\}$.

accuracy edge

Suppose $\|Y - f^*(X)\|_{L_2} \leq \sigma$ for a known constant $\sigma > 0$.
Introduce the “complexity”

$$r^* = \max\{\lambda_{\mathbb{Q}}(c_1, c_2), \lambda_{\mathbb{M}}(c_1/\sigma, c_2), r_E(c_1), \tilde{r}_{\mathbb{M}}(c_1, \sigma)\} .$$

Mendelson (2016) proved that r^* is an upper bound for the accuracy edge (under a “small-ball” assumption).

linear regression—an example

Let $\mathcal{F} = \{\langle \mathbf{t}, \cdot \rangle : \mathbf{t} \in \mathbb{R}^d\}$ be the class of linear functionals.

Let \mathbf{X} be an isotropic random vector in \mathbb{R}^d such that
 $\|\langle \mathbf{X}, \mathbf{t} \rangle\|_{L_4} \leq L \|\langle \mathbf{X}, \mathbf{t} \rangle\|_{L_2}.$

Suppose $\mathbf{Y} = \langle \mathbf{t}_0, \mathbf{X} \rangle + \mathbf{W}$ for some $\mathbf{t}_0 \in \mathbb{R}^d$ and symmetric independent noise \mathbf{W} with variance σ^2 .

linear regression

Given n independent samples $(\mathbf{X}_i, \mathbf{Y}_i)$, least-squares regression (ERM) finds $\hat{\mathbf{t}}_n$ such that

$$\|\hat{\mathbf{t}}_n - \mathbf{t}\| \leq c \frac{\sigma}{\delta} \sqrt{\frac{d}{n}}$$

with probability $1 - \delta - e^{-cd}$.

Note the weak accuracy/confidence tradeoff.

Lecué and Mendelson (2016) show that this is essentially optimal.

However, if everything is sub-Gaussian, one has

$$\|\hat{\mathbf{t}}_n - \mathbf{t}\| \leq c\sigma \sqrt{\frac{d}{n}}$$

with probability $1 - e^{-cd}$.

We introduce a procedure that achieves the same performance as sub-Gaussian ERM but under the general fourth-moment condition.

median-of-means tournament

A natural idea is to replace ERM by minimization of the median-of-means estimate of the risk $\mathbb{E}(f(\mathbf{X}) - Y)^2$.

Difficult to analyze—may be suboptimal.

median-of-means tournament

A natural idea is to replace ERM by minimization of the median-of-means estimate of the risk $\mathbb{E}(\mathbf{f}(\mathbf{X}) - \mathbf{Y})^2$.

Difficult to analyze—may be suboptimal.

Instead, we run a **median-of-means tournament**.

The idea is that, based on a median-of-means estimate of the **difference**

$$\mathbb{E}(\mathbf{f}(\mathbf{X}) - \mathbf{Y})^2 - \mathbb{E}(\mathbf{h}(\mathbf{X}) - \mathbf{Y})^2 ,$$

we can have a good guess if **f** or **h** has a smaller risk.

median-of-means tournament

To make the idea work, we design a (two- or) three-step procedure.

Each step uses an independent sample so before starting we split the data into (two or) three equal parts.

The procedure has a parameter $r > 0$, the desired accuracy level.

The main steps of the procedure are:

- Distance referee
- Elimination phase
- Champions league

step 1: the distance referee

For each pair $f, h \in \mathcal{F}$, one may use define a median-of-means estimate $\Phi_n(f, h)$ using $(|f(X_i) - h(X_i)|)_{i=1}^n$ such that, with “high probability”, for all $\Phi_n(f, h)$,

$$\text{if } \Phi_n(f, h) \geq \beta r \text{ then } \|f - h\|_{L_2} \geq r$$

and

$$\text{if } \Phi_n(f, h) < \beta r \text{ then } \|f - h\|_{L_2} < \alpha r$$

for some constants α, β .

Matches are only allowed between $f, h \in \mathcal{F}$ if $\Phi_n(f, h) \geq \beta r$.

step 2: elimination phase

For any pair $f, h \in \mathcal{F}$, if the distance referee allows a match, calculate the median-of-means estimate based on the samples

$$(f(X_i) - Y_i)^2 - (h(X_i) - Y_i)^2.$$

if the estimate is negative, f wins the match otherwise h wins.

$f \in \mathcal{F}$ is a champion if it wins all its matches. Let \mathcal{H} be the set of all champions.

If one only cares about the mean squared error $\|\hat{f}_n - f^*\|_{L_2}$, then one may select any champion $\hat{f}_n \in \mathcal{H}$.

One may show that, with “high probability”, \mathcal{H} contains f^* and possibly other functions within distance $O(r)$ of f^* .

If the excess risk also matters, all champions in \mathcal{H} advance to the Champions League for the playoffs.

step 3: Champions League

To select a champion with a small excess risk, we use the simple fact that, for any $f \in \mathcal{F}$,

$$\begin{aligned} \mathbb{E}(f(X) - Y)^2 - \mathbb{E}(f^*(X) - Y)^2 \\ \leq -2\mathbb{E}(f^*(X) - f(X))(f(X) - Y) . \end{aligned}$$

The Champions League winner is selected based on median-of-means estimates of $\mathbb{E}(h(X) - f(X))(f(X) - Y)$ for all pairs $f, h \in \mathcal{F}$.

result

Suppose that \mathcal{F} is a convex class of functions and

- for every $f, h \in \mathcal{F}$, $\|f - h\|_{L_4} \leq L\|f - h\|_{L_2}$;
- for every $f \in \mathcal{F}$, $\|f - Y\|_{L_4} \leq L\|f - Y\|_{L_2}$;

Then the median-of-means tournament achieves an essentially optimal accuracy/confidence tradeoff.

For any $r > r^*$, with probability at least

$$1 - \exp(-c_0 n \min\{1, \sigma^{-2} r^2\}) ,$$

$$\|\hat{f} - f^*\|_{L_2} \leq cr$$

and

$$\mathbb{E}((\hat{f}(X) - Y)^2 | \mathcal{D}_n) \leq \mathbb{E}(f^*(X) - Y)^2 + (cr)^2 .$$

linear regression

Recall the example $\mathcal{F} = \{\langle t, \cdot \rangle : t \in \mathbb{R}^d\}$ with \mathbf{X} isotropic such that $\|\langle \mathbf{X}, t \rangle\|_{L_4} \leq L \|\langle \mathbf{X}, t \rangle\|_{L_2}$ and $\mathbf{Y} = \langle t_0, \mathbf{X} \rangle + \mathbf{W}$.

We obtain

$$\|\hat{t}_n - t\| \leq c\sigma \sqrt{\frac{d}{n}}$$

with probability $1 - e^{-cd}$ and also

$$\mathbb{E}((\hat{f}(\mathbf{X}) - \mathbf{Y})^2 | \mathcal{D}_n) - \mathbb{E}(f^*(\mathbf{X}) - \mathbf{Y})^2 \leq c\sigma^2 \frac{d}{n} .$$

algorithmic challenge

Find an algorithmically efficient version of the median-of-means tournament.

references

- G. Lugosi and S. Mendelson.
Sub-Gaussian estimators of the mean of a random vector.
submitted, 2017.
- G. Lugosi and S. Mendelson.
Risk minimization by median-of-means tournaments.
submitted, 2016.
- E. Joly, and G. Lugosi, and R. Imbuzeiro Oliveira.
On the estimation of the mean of a random vector.
Electronic Journal of Statistics, 2017.
- L. Devroye, M. Lerasle, G. Lugosi, and R. Imbuzeiro Oliveira.
Sub-Gaussian mean estimators.
Annals of Statistics, 2016.

references

C. Brownlees, E. Joly, and G. Lugosi.

Empirical risk minimization for heavy-tailed losses.

Annals of Statistics, 43:2507–2536, 2015.

E. Joly, and G. Lugosi.

Robust estimation of U-statistics.

Stochastic Processes and their Applications, to appear, 2015.

S. Bubeck, N. Cesa-Bianchi, and G. Lugosi.

Bandits with heavy tail.

IEEE Transactions on Information Theory, 59:7711–7717, 2013.