

# Recent Advances in Algorithmic Heavy-Tailed Statistics

Sam Hopkins, UC Berkeley (→ MIT)

Based on joint work with Yeshwanth Cherapanamjeri, Tarun Kathuria, Prasad Raghavendra, and Nilesh Tripuraneni.

# Measuring success probability in estimation (confidence intervals)

Given  $X_1, \dots, X_n \sim P_\theta$ , find  $\hat{\theta}$  s.t.  $\|\theta - \hat{\theta}\| \leq r(n, d, \delta)$   
With prob.  $\geq 1 - \delta$

no. samples  $\nearrow$  "rate"  
ambient dimension  $\uparrow$   
failure rate  $\nwarrow$

For  $P_\theta \in \mathcal{P}$   $\longleftarrow$  class of distributions

Our game: large class  $\mathcal{P}$  (e.g. all distributions with  $O(1)$  bdd. moments), but similar guarantees as if  $\mathcal{P}$  contains Gaussians.

Example: Estimating the mean in  $\ell_2$

$$X_1, \dots, X_n \sim X, \text{Cov}(X) \preceq I$$

$$\mathbb{E} \left\| \frac{1}{n} \sum X_i - \mathbb{E}X \right\|^2 \leq e^{-t^2 n}$$

$$\text{Gaussian: } \left\| \frac{1}{n} \sum X_i - \mathbb{E}X \right\| \leq \sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}} \quad \text{w.p. } 1-\delta$$

---

Drop the Gaussian assumption?

Is there an estimator  $\hat{\mu}_\delta(x_1, \dots, x_n)$  s.t.

$$\|\hat{\mu}_\delta - \mu\| \leq O\left(\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}}\right) \quad \text{w.p. } 1-\delta?$$

## Example: Estimating the mean in $\ell_2$

$d=1$ : Median of means, truncated mean, ...

Lugosi-Mendelson '18: estimator based on high-dimensional quantiles/median of means (exp. time)  $\text{poly}(n, d, \frac{1}{\delta})$

H. '18: algorithmic approach to high-dimensional quantiles

Theorem: for every  $\delta > 2^{-n}$  exists  $\hat{\mu}_\delta(x_1 \dots x_n)$  s.t.

$$\|\hat{\mu}_\delta - \mu\| \leq O\left(\sqrt{\frac{d}{n}} + \sqrt{\frac{\log(1/\delta)}{n}}\right) \text{ w.p. } 1-\delta, \text{ poly}(n, d, \log(1/\delta))$$

time

# Agenda

1. Survey of recent (algorithmic) developments
2. Heavy-tailed covariance estimation

- Algorithm of [Cherapanamjeri - H. - Kathuria - Raghavendra - Tripuraneni]

- 2 key techniques for concentration + convex prog.

- bounded differences

- SoS Bernstein

$$\|\Sigma M_i\|_{op} = \max_{\substack{Y \neq 0 \\ \|Y\|_2 = 1}} \langle Y, \Sigma M_i \rangle$$

$$\|\Sigma M_i\|_{SoS} = \max_{Y \in SoS} \langle Y, \Sigma M_i \rangle$$

# Matching Gaussian Confidence w/ only $O(1)$ moments

- mean estimation,  $\ell_2$  - median of means, trimmed mean
- mean estimation,  $\|\cdot\|$
- COVARIANCE estimation
- regression
- regularized regression/sparse recovery
- . . .

[Lugosi-Mendelson '19, '19, '19, '20, '20, Mendelson-Zhivotovskiy '20, Minsker '15, '18, '20, Hsu-Sabato '16, Lerasale-Oliveira '12, Joly-Lugosi-Oliveira '17, ... ]

# Matching Gaussian Confidence w/ only $O(1)$ moments

in polynomial time

Mean estimation,  $l_2$

authors	time
---------	------

H.

$O(nd)^{28}$

(Cherapanamjeri -  
Flammann - Bartlett)

$\tilde{O}(n^{35} + n^2d)$

Depersin - Lecué,  
Lei-Luh-Venkat-Zhang

$\tilde{O}(n^2d)$

Covariance,  $\|\cdot\|_p$

authors	Suboptimal by
---------	---------------

Minsker

$\sqrt{\log 1/s}$

CHKRT

$(\log 1/s)^{1/4}$

Similar for linear regression

(Hsu-Sabato)

# Connections to Robust Statistics



- related techniques – new notions of high-dimensional quantiles & algorithms to compute them
- Simultaneously robust and high-confidence algos

[Depersin-Lecue '19, Diakonikolas-Kane-Pensia '20, H.-Li-Zhang '20,...]



# Agenda

- ~~1. Survey of recent (algorithmic) developments~~
2. Heavy-tailed covariance estimation
  - Algorithm of [Cherapanamjeri - H. - Kathuria - Raghavendra - Tripuraneni]
  - 2 Key techniques for concentration + convex prog.
    - bounded differences
    - SoS Bernstein

Setup:  $X_1, \dots, X_n \sim X$  on  $\mathbb{R}^d$ , covariance  $\Sigma \in \mathbb{R}^{d \times d}$   
 for simplicity,  $\text{Tr} \Sigma = O(d)$ ,  $\|\Sigma\|_{\text{op}} = O(1)$

Theorem: Can find  $\hat{\Sigma}_S$  s.t.  $\|\hat{\Sigma}_S - \Sigma\|_{\text{op}} \leq \tilde{O}\left(\sqrt{\frac{d}{n}} \cdot (\log' 1/\delta)^{1/4} + \sqrt{\frac{\log' 1/\delta}{n}}\right)$   
 w.p.  $1 - \delta$  in time  $\text{poly}(n, d, \log' 1/\delta)$  if  $X$  is **Certifiably (2, 8)-hypercontractive**.

• Sub-optimal by  $(\log' 1/\delta)^{1/4}$  - can remove in:

- exp. time [Mendelson-Zhivotovskiy]
- $\text{poly}(n, d, \frac{1}{\delta})$  time [CHKRT]

• best previous:  $\tilde{O}\left(\sqrt{\frac{d \log' 1/\delta}{n}}\right)$  [Minsker]

only require **(2, 4)-hypercontractivity**

Theorem: Can find  $\hat{\Sigma}_S$  s.t.  $\|\hat{\Sigma}_S - \Sigma\|_{\text{op}} \leq \tilde{O}\left(\sqrt{\frac{d}{n}} \cdot (\log 1/\delta)^{1/4} + \sqrt{\frac{\log 1/\delta}{n}}\right)$   
w.p.  $1-\delta$  in time  $\text{poly}(n, d, \log 1/\delta)$  if  $X$  is certifiably  $(2, 8)$ -hypercontractive.

$(2, 8)$ -hypercontractivity:  $\mathbb{E}\langle X_{1,v} \rangle^8 \leq O(\mathbb{E}\langle X_{1,v}^2 \rangle)^4$

Certifiable:  $O(\mathbb{E}\langle X_{1,v}^2 \rangle)^4 - \mathbb{E}\langle X_{1,v} \rangle^8 = \sum p_i \langle v \rangle^2$

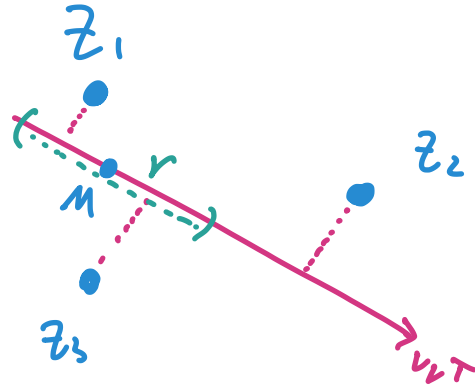
- products of heavy-tailed distns
- mixtures thereof
- affine transformations thereof
- ...



# The High-Dimensional Median-of-Means Paradigm

"Spectral  $r$ -center":

$\forall v, |\langle z_i, vv^T \rangle - \langle M, vv^T \rangle| \leq r$   
for 60% of  $z_i$ 's.



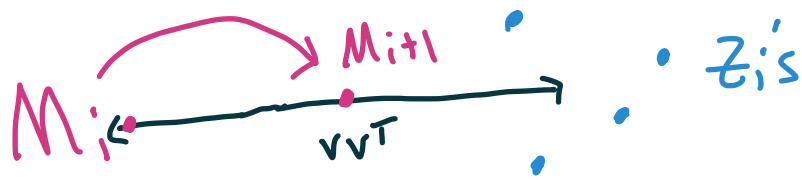
- $\Sigma$  is a spectral  $\sqrt{\frac{d}{n}} + \sqrt{\frac{\log 1/\delta}{n}}$  center of  $z_1, \dots, z_{\log 1/\delta}$  w.p.  $1 - \delta$   
[Lugosi-Mendelson, Mendelson-Zhiva.]

- $M, M'$  spectral  $r$ -centers  $\rightarrow \|M - M'\|_{op} \leq 2r$

Just need to find a spectral  $r$ -center!

Our Strategy:  $M_1 \rightarrow M_2 \rightarrow \dots \rightarrow M$

- Attempt to certify that  $M_i$  is a spectral  $r$  center
- Success  $\rightarrow$  output  $M_i$
- Failure  $\rightarrow$  witness  $v$ , update  $M_{i+1} = M_i \pm vv^T$



[CFB '19, CHKRT '20]

If GIVEN  $\Sigma$  and  $z_1, \dots, z_{\log 1/s}$ ,

can you check that

$$|\langle z_i, vv^T \rangle - \langle \Sigma, vv^T \rangle| \leq r$$

for 60% of  $z_i$ 's and all unit  $v$ ?

# Optimization Approach

Variables  $b_1, \dots, b_{\log_2 1/s}, v_1, \dots, v_d$

$$\max_{b, v} \sum_{i=1}^{\log_2 1/s} b_i$$

$$\text{s.t. } b_i \in \{0, 1\}, \|v\| = 1,$$

$$b_i (\langle z_i, v v^T \rangle - \langle M, v v^T \rangle) \geq b_i \cdot r$$

If GIVEN  $\Sigma$  and  $z_1, \dots, z_{\log_2 1/s}$ ,  
can you check that

$$|\langle z_i, v v^T \rangle - \langle \Sigma, v v^T \rangle| \leq r$$

for 60% of  $z_i$ 's and all unit  $v$ ?



# Optimization Approach

Variables  $b_1, \dots, b_{\log_2 s}, v_1, \dots, v_d = v$

$$\max_{b, v} \sum_{i=1}^{\log_2 s} b_i$$

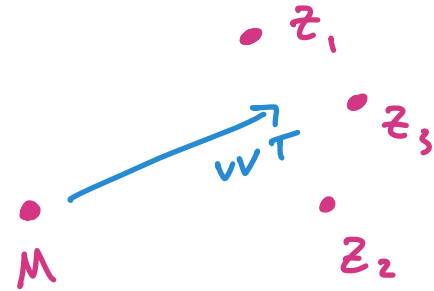
$$\text{s.t. } b_i \in \{0, 1\}, \|v\| = 1,$$

$$b_i (\langle z_i, v v^T \rangle - \langle M, v v^T \rangle) \geq b_i \cdot r$$

If GIVEN  $\Sigma$  and  $z_1, \dots, z_{\log_2 s}$ ,  
can you check that

$$|\langle z_i, v v^T \rangle - \langle \Sigma, v v^T \rangle| \leq r$$

for 60% of  $z_i$ 's and all unit  $v$ ?



# Optimization Approach

Variables  $b_1, \dots, b_{\log^2 s}, v_1, \dots, v_d$

$$\max_{b, v} \sum_{i=1}^{\log^2 s} b_i$$

$$\text{s.t. } b_i \in \{0, 1\}, \|v\| = 1,$$

$$b_i (\langle z_i, v v^T \rangle - \langle M, v v^T \rangle) \geq b_i \cdot r$$

- Convex relaxation:  $\leq 0.6 \log^2 s$ , for  $r = (\log^2 s)^{1/4} \sqrt{\frac{d}{n}} + \sqrt{\frac{\log^2 s}{n}}$ , w.p.  $1 - \delta$

If GIVEN  $\Sigma$  and  $z_1, \dots, z_{\log^2 s}$ ,  
can you check that

$$|\langle z_i, v v^T \rangle - \langle \Sigma, v v^T \rangle| \leq r$$

for 60% of  $z_i$ 's and all unit  $v$ ?

- Know:  $\leq 0.6 \cdot \log^2 s$  for  $r = \sqrt{\frac{d}{n}} + \sqrt{\frac{\log^2 s}{n}}$ , w.p.  $1 - \delta$

# Agenda

- ~~1. Survey of recent (algorithmic) developments~~
2. Heavy-tailed covariance estimation
  - ~~Algorithm of [Cherapanamjeri - H. - Kathuria - Raghavendra - Tripuraneni]~~
  - 2 key techniques for concentration + convex prog.
    - bounded differences
    - SoS Bernstein

# Optimization Approach

Variables  $b_1, \dots, b_{\log^2/s}, v_1, \dots, v_d$

$$\max_{b, v} \sum_{i=1}^{\log^2/s} b_i$$

$$\text{s.t. } b_i \in \{0, 1\}, \|v\| = 1,$$

$$b_i (\langle z_i, v v^T \rangle - \langle \mathbb{E} z_i, v v^T \rangle) \geq b_i \cdot r$$

If GIVEN  $\Sigma$  and  $z_1, \dots, z_{\log^2/s}$ ,  
Can you check that

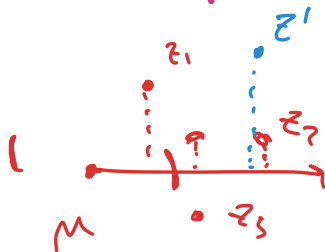
$$|\langle z_i, v v^T \rangle - \langle \mu, v v^T \rangle| \leq r$$

for 60% of  $z_i$ 's and all unit  $v$ ?

Any "reasonable" convex relaxation will satisfy bounded diffs.

$$|\widetilde{\text{OPT}}(\underline{z_1, \dots, z_{\log^2/s}}) - \widetilde{\text{OPT}}(\underline{z'_1, \dots, z'_{\log^2/s}})| \leq 1$$

$$B_i \sim b_i \quad 0 \leq B_i \leq 1$$



New Goal: Convex relaxation of

$$z_i = \mathbb{E}_{j \sim B_i} x_j x_j^T$$

$$\max_{b, v} \sum_{i=1}^{\log' / s} b_i$$

$$\text{s.t. } b_i \in \{0, 1\}, \|v\| = 1,$$

$$b_i (\langle z_i, v v^T \rangle - \langle \Sigma, v v^T \rangle) \geq b_i \cdot r$$

$$\text{s.t. } \mathbb{E}_{z_1, \dots, z_{\log' / s}} \text{OPT} \leq 0.6 \log' / s, \text{ assuming only}$$

Certifiable  $(2, \delta)$ -hypercontractivity

We use Sum of Squares Semidefinite Program

hierarchy of convex relaxations for any polynomial optimization problem

$$\begin{aligned} \max_{b, v} \quad & \sum_{i=1}^{\log 1/s} b_i \\ \text{st.} \quad & b_i(1-b_i) = 0 = \sum v_i^2 \quad \text{proxy variables for } b_i, b_i b_j, \underbrace{v_i v_j b_k}_{B_i \dots} \\ & b_i \in \{0, 1\}, \quad \|v\|^2 = 1 \\ & \underbrace{b_i \langle z_i, v v^T \rangle - \langle a, v v^T \rangle}_{B_i} \geq b_i \cdot r \end{aligned}$$

①

$$\max_{b, v} \sum_{i=1}^{\log^2 \delta} b_i$$

$$\text{st. } b_i \in \{0, 1\}, \|v\| = 1,$$

$$b_i (\langle z_i, v v^T \rangle - \langle \Sigma, v v^T \rangle) \geq b_i \cdot r$$

$$\leq \max_{b, v} \frac{1}{r} \sum b_i (\langle z_i, v v^T \rangle - \langle \Sigma, v v^T \rangle)$$

$$\leq \max_{\|v\|=1} \frac{1}{r} \sqrt{\log^2 \delta} \cdot \left( \sum (\langle z_i - \Sigma, v v^T \rangle^2)^{\frac{1}{2}} \right)$$

Suffices to bound

$$\mathbb{E}_{z_1, \dots, z_{\log^2 \delta}} \text{SoS} \left( \max_{\|v\|=1} \sum_{i=1}^{\log^2 \delta} \langle z_i - \Sigma, v v^T \rangle^2 \right)$$

Suffices to bound

$$\mathbb{E} \text{SoS} \left( \max_{\|v\|=1} \sum_{i=1}^{\log' s} \langle z_i - \Sigma, vv^T \rangle^2 \right)$$

$z_1 \dots z_{\log' s}$

Sum of i.i.d. polynomials in  $v$



Suffices to bound

$$\mathbb{E} \text{SoS} \left( \max_{\|v\|=1} \sum_{i=1}^{\log'1s} \langle z_i - \Sigma, vv^T \rangle^2 \right)$$

Sum of i.i.d. polynomials in  $v$

$$\max_{Y \in \mathbb{R}^{d^2 \times d^2}} \int \langle (z_i - \Sigma) \otimes (z_i - \Sigma), Y \rangle$$

$Y \in \text{SoS}$

$\text{Tr} Y = 1, Y \succeq 0$

Suffices to bound

$$\mathbb{E} \text{SoS} \left( \max_{\|v\|=1} \sum_{i=1}^{\log' 1/s} \langle z_i - \Sigma, v v^T \rangle^2 \right)$$

Sum of i.i.d. polynomials in  $v$

$$\mathbb{E} \max_{Y \in \mathbb{R}^{d^2 \times d^2}} \left[ \langle (z_i - \Sigma) \otimes (z_i - \Sigma), Y \rangle \right]$$

$\{ Y \in \text{SoS} \}$

$$= \mathbb{E} \left\| \sum (z_i - \Sigma) \otimes (z_i - \Sigma) \right\|_{\text{SoS}}$$

$$\textcircled{2} \quad \mathbb{E} \left\| \sum (z_i - \Sigma) \otimes (z_i - \Sigma) \right\|_{\text{SoS}}$$

SoS Bernstein: control via  $\mathbb{E} \left\| \sum (z_i - \Sigma)^2 \otimes (z_i - \Sigma)^2 \right\|_{\text{SoS}}$

degree 8 polynomial in  $X$

$\Rightarrow$  control via (2.8) certifiable

hypercontractivity:

$$\mathbb{E} \| X^{\otimes 8} \|_{\text{SoS}} \leq O(1)$$

Yields covariance estimation algo. with  $r = \tilde{O}\left((\log'ls)^{1/4} \sqrt{\frac{d}{n}} + \sqrt{\frac{\log'ls}{n}}\right)$

Open: remove  $(\log'ls)^{1/4}$  ?

Open: optimal heavy-tailed linear regression?

Applications of

- bdd. difs for convex programs?
- SoS Bernstein?

E  $\langle x, v \rangle^8$

T H A N K S