

Course: High Dimensional Probability:  
applications to statistical learning  
and data science.

Text(s): "High-Dimensional Probability"  
Roman Vershynin  
"High-Dimensional Statistics"  
Martin J. Wainwright

Geography: Let  $n$  be sample size  
 $d$  is the dimension of data

- Classical Asymptotics

$n \rightarrow \infty$ ,  $d$  fixed

Ex: Laws of large numbers  
Central limit theorem

- High dimensional Asymptotics

$(n, d) \rightarrow \infty$  while  $\psi(n, d) \rightarrow \alpha < \infty$

where  $\psi$  is a "scaling function"

E.g.:  $\psi(n, d) = \frac{d}{n}$

- Non-asymptotic Bounds  $\leftarrow$  This class  
 $(n, d)$  is fixed. Probability of successful  
estimation/inference depends on  $(n, d)$

See (W:1) for examples of the three viewpoints.

Lecture 1: Review (V:1)

Let  $X$  be a random variable on probability space

Expectation and Variance

$$EX, \text{Var}(X) = E(X - EX)^2 = E[X^2] - (EX)^2$$

Moment Generating Function:

$$M_X(t) = Ee^{tX}, t \in \mathbb{R}$$

$L^p$ -norm  $\|X\|_p = (E|X|^p)^{1/p}$ ,  $p \in (0, \infty)$

Banach space:

$$L^p = \{X : \|X\|_p < \infty\}$$

Remark:  $L^2$  is a Hilbert space

$$\langle X, Y \rangle_2 = E[XY], \quad \|X\|_2 = \sqrt{\langle X, X \rangle} = \sqrt{E[X^2]}$$

Then

$$\|X - EX\|_2 = \sqrt{\text{Var}(X)} \quad \text{and}$$

$$\begin{aligned} \text{cov}(X, Y) &:= E(X - EX)(Y - EY) \\ &= \langle X - EX, Y - EY \rangle_2 \end{aligned}$$

## Limit Theorems

Thm: (Strong Law of Large Numbers)

Let  $X_1, X_2, \dots$  be i.i.d. random variables with  $\mathbb{E}X_i = \mu$ . Then  $S_n = \sum_{i=1}^n X_i$  satisfies

$$\frac{1}{n} \cdot S_n \rightarrow \mu \text{ almost surely}$$

$$\left[ \lim_{n \rightarrow \infty} \frac{S_n}{n} = \mu \text{ w.p. } 1 \right]$$

Intuition:  $\text{Var}(S_n) = \frac{\text{Var}(X_i)}{n}$

Thm: (Central Limit Theorem)

Let  $X_1, X_2, \dots$  be i.i.d. with

$$\mathbb{E}X_i = \mu, \quad \text{Var} X_i = \sigma^2$$

Define  $Z_n = \frac{S_n - \mathbb{E}S_n}{\sqrt{\text{Var}(S_n)}} = \frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n (X_i - \mu)$

Then  $Z_n \rightarrow N(0, 1)$  in distribution

$$\left[ \mathbb{P}[Z_n \geq t] \rightarrow \frac{1}{\sqrt{2\pi}} \int_t^{\infty} e^{-x^2/2} dx \right]$$

# Chapter 1

## Concentration Inequalities

- Chernoff Bound
- Sub-Gaussian RV (Hoeffding)
- Sub-Exponential RV (Bernstein)
- Johnson-Lindenstrauss
- McDiarmid inequality
- Robust Mean Estimation

## Concentration Inequalities

Goal: Bounds of the form

$$P\left[\frac{S_n}{n} > \mu + t\right] \leq \text{small}(n, t)$$

Prop: Let  $Z \sim N(\mu, \sigma)$ . Then

$$P[Z > \mu + t] \leq e^{-t^2/2\sigma^2} \quad \forall t > 0$$

So would hope  $P\left[\frac{S_n}{n} > \mu + t\right] \leq C e^{-\frac{nt^2}{2\sigma^2}}$

**Bad news from CLT:**

$$\sup_t \left| P[Z_n > t] - P[Z > t] \right| \left\{ \begin{array}{l} \text{can be } \Omega\left(\frac{1}{\sqrt{n}}\right) \\ \text{even for Bernoulli} \end{array} \right.$$

where  $Z \sim N(0, 1)$

[See (V: 2.1) for example]

Chernoff Method: {Sub-Gaussian} RV's  
{Sub-Exponential}

Lemma (Markov) For any non-negative  $X$  and  $t > 0$ , we have

$$P[X \geq t] \leq \frac{E[X]}{t}$$

pf:  $E[X] = E[X \mathbb{1}_{\{X \geq t\}}] + E[X \mathbb{1}_{\{X < t\}}]$

$$\geq t E[\mathbb{1}_{\{X \geq t\}}] + 0 = t P[X \geq t]$$

Cor: (Chebychev) For any random variable  $X$  with  $\text{var}(X) \leq \sigma^2 < \infty$ , we have

$$P[|X - E[X]| \geq t] \leq \frac{\sigma^2}{t^2} \quad \forall t > 0$$

More generally, suppose  $\mu = E[X] < \infty$ . Then for all  $\lambda > 0$ , we have

$$P[X - \mu \geq t] = P[e^{\lambda(X-\mu)} \geq e^{\lambda t}] \stackrel{\text{Markov}}{\leq} \frac{E[e^{\lambda(X-\mu)}]}{e^{\lambda t}}$$

$$\Rightarrow \log P[X - \mu \geq t] \leq \inf_{\lambda > 0} \left\{ \log E[e^{\lambda(X-\mu)}] - \lambda t \right\} \\ = - \sup_{\lambda > 0} \left\{ \lambda t - \log E[e^{\lambda(X-\mu)}] \right\}$$

Define for any function  $\psi: \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ , the Fenchel conjugate

$$\psi^*(t) = \sup_{\lambda} \{t\lambda - \psi(\lambda)\}$$

Let's look at the main example

$$\psi_x(\lambda) = \log \mathbb{E} e^{\lambda(x-\mu)}$$

For all  $\lambda \in \mathbb{R}$ , observe

$$\psi_x(\lambda) = \log \mathbb{E} e^{\lambda(x-\mu)} \geq \mathbb{E} \log e^{\lambda(x-\mu)} = 0$$

So when  $\lambda < 0$  and  $t > 0$ , we have

$$t\lambda - \psi(\lambda) \leq 0 = 0 \cdot t - \psi(0).$$

Therefore for  $t \geq 0$ , equality holds:

$$\psi_x^*(t) = \sup_{\lambda \geq 0} \{t\lambda - \psi(\lambda)\}$$

We therefore arrive at the bound:

Chernoff Bound:

$$\mathbb{P}[X - \mu \geq t] \leq \exp(-\psi_X^*(t))$$

where  $\psi_X(\lambda) := \log(\mathbb{E} e^{\lambda(X-\mu)})$ .

Ex: Let  $X \sim N(\mu, \sigma^2)$ . Then

$$\mathbb{E} e^{\lambda(X-\mu)} = e^{\frac{\sigma^2 \lambda^2}{2}} \quad \forall \lambda \in \mathbb{R}.$$

$$\text{So } \mathbb{P}[X \geq \mu + t] \leq e^{-\frac{t^2}{2\sigma^2}} \quad \forall t > 0$$

Defn:  $X$  with mean  $\mu = \mathbb{E}X$  is sub-Gaussian with parameter  $\sigma > 0$  if

$$\mathbb{E} e^{\lambda(X-\mu)} \leq e^{\frac{\sigma^2 \lambda^2}{2}} \quad \forall \lambda \in \mathbb{R}$$

If  $X$  is sub-Gaussian, so is  $-X$ .

$$\Rightarrow \mathbb{P}[|X - \mu| \geq t] \leq 2 e^{-\frac{t^2}{2\sigma^2}}$$



Ex: (Rademacher)

Let  $\varepsilon$  be a rademacher RV:

$$P[\varepsilon = 1] = P[\varepsilon = -1] = \frac{1}{2}.$$

$$\mathbb{E}[e^{\lambda \varepsilon}] = \frac{1}{2}(e^{-\lambda} + e^{\lambda}) = \cosh(\lambda)$$

Exercise:  $\cosh(x) \leq \exp(\frac{x^2}{2}) \quad \forall x \in \mathbb{R}$

So  $\varepsilon$  is sub-Gaussian with  $\sigma=1$ .

Ex: (Bounded RV)

Suppose  $X$  is supported on  $[a, b]$ .

Then Jensen's inequality implies

$$\mathbb{E} e^{\lambda(X - \mathbb{E}X)} \leq \mathbb{E} e^{\lambda(X - X')}$$

where  $X'$  is independent copy of  $X$ .

Observe  $X - X' \sim \varepsilon(X - X')$

$$\text{So } \mathbb{E}_{X, X'} e^{\lambda(X - X')} = \mathbb{E}_{X, X'} \mathbb{E}_{\varepsilon} e^{\varepsilon \cdot \lambda(X - X')}$$

$$\leq \mathbb{E}_{X, X'} \exp(\lambda^2 (X - X')^2 / 2)$$

$$\leq \exp(\lambda^2 (b-a)^2 / 2)$$

A slightly more involved argument shows that  $X$  is sub-Gaussian with parameter  $\frac{b-a}{2}$

$$\Rightarrow \mathbb{P}[|X-\mu| \geq t] \leq 2 \exp\left(-\frac{2t^2}{(b-a)^2}\right)$$

Lemma: [Sum Rule]

$X_i$  are independent  $\sigma_i$ -sub-Gaussian  $\forall i=1, \dots, n$   $\Rightarrow \sum_{i=1}^n X_i$  is  $\|\sigma\|_2$ -sub-Gaussian

Cor: (Hoeffding)

Suppose  $X_1, \dots, X_n$  are independent with  $\mathbb{E}X_i = \mu_i$  and are  $\sigma_i$ -sub-Gaussian

Then

$$\mathbb{P}\left[\sum_{i=1}^n (X_i - \mu_i) \geq t\right] \leq \exp\left\{-\frac{t^2}{2\|\sigma\|_2^2}\right\}$$

$\Rightarrow$  If  $\mu_i = \mu$ ,  $\sigma_i = \sigma$ , then

$$\mathbb{P}\left[\frac{1}{n} \sum_{i=1}^n X_i \geq \mu + t\right] \leq \exp\left\{-\frac{nt^2}{2\sigma^2}\right\}$$

Subexponential RV:

Ex: Let  $Z \sim N(0,1)$ . Let's compute

$$\begin{aligned}\mathbb{E}[e^{\lambda(Z^2-1)}] &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{\lambda(x^2-1)} e^{-x^2/2} dx \\ &= \frac{e^{-\lambda}}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-(1-2\lambda)x^2/2} dx \\ &= \begin{cases} \frac{e^{-\lambda}}{\sqrt{1-2\lambda}} & \text{if } \lambda \leq \frac{1}{2} \\ +\infty & \text{if } \lambda > \frac{1}{2} \end{cases}\end{aligned}$$

Defn:  $X$  with mean  $\mu = \mathbb{E}X$  is subexponential with parameters  $(\sigma, \alpha)$  if  $\mathbb{E}[e^{\lambda(X-\mu)}] \leq e^{\frac{\sigma^2 \lambda^2}{2}} \forall |\lambda| \leq \frac{1}{\alpha}$

Back to example  $Z \sim N(0,1)$

$$\mathbb{E}[e^{\lambda(Z^2-1)}] \leq \frac{e^{-\lambda}}{\sqrt{1-2\lambda}} \leq e^{\frac{4\lambda^2}{2}} \quad \forall |\lambda| < \frac{1}{4}$$

So  $Z^2$  is  $(2, 4)$ -subexponential.

Thm (Subexponential tail bound)

Let  $X$  be subexponential with  $(\sigma, \alpha)$ .

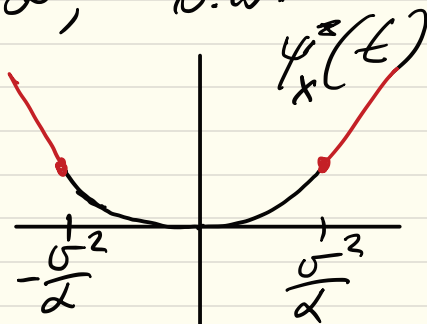
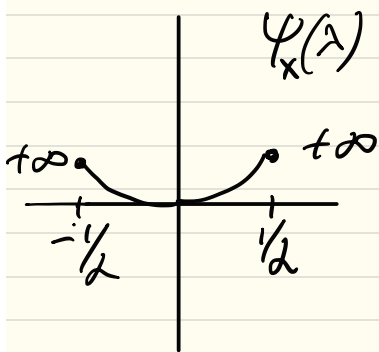
Then

$$P[X - \mu \geq t] \leq \begin{cases} e^{-\frac{t^2}{2\sigma^2}}, & \text{if } |t| \leq \frac{\sigma^2}{2} \\ e^{-\frac{t}{2\alpha}}, & \text{o.w.} \end{cases}$$

pf: Back to Chernoff

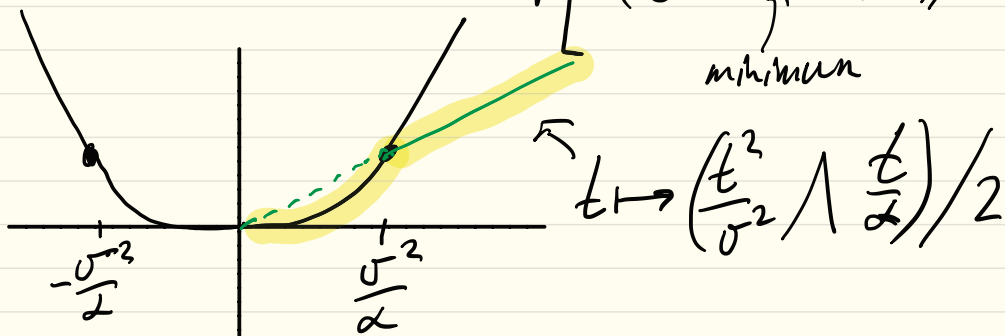
$$\log P[X - \mu \geq t] \leq -\Psi_X^*(t)$$

$$\text{where } \Psi_X(\lambda) = \log E[e^{\lambda(X - \mu)}] \\ = \begin{cases} \frac{\sigma^2 \lambda^2}{2}, & \text{if } |\lambda| \leq \frac{1}{\sigma} \\ +\infty, & \text{o.w.} \end{cases}$$



Thm: (Bernstein) Let  $X$  be subexponential with parameter  $(\sigma, \alpha)$  and mean  $\mu = \mathbb{E}X$ . Then

$$\mathbb{P}[|X - \mu| \geq t] \leq 2 \exp\left[-\left(\frac{t^2}{\sigma^2} \wedge \frac{t}{\alpha}\right) / 2\right]$$



Lemma: [Sum Rule]

$X_i$  are  $(\sigma_i, \alpha_i)$ -subexp  $\Rightarrow \sum_{i=1}^n X_i$  is  $(\|\sigma\|_2, \|\alpha\|_\infty)$ -subexp

Thm (Bernstein for sums) Let  $X_1, \dots, X_n$  be independent subexponential with parameters  $(\sigma_i, \alpha_i)$ , and with mean  $\mu_i = \mathbb{E}X_i$ .

Then

$$\mathbb{P}\left[\left|\sum_{i=1}^n (X_i - \mu_i)\right| \geq t\right] \leq 2 \exp\left[-\frac{1}{2} \left(\frac{t^2}{\|\sigma\|_2^2} \wedge \frac{t}{\|\alpha\|_\infty}\right)\right]$$

Thm: (Improved Bernstein for bounded RV)

Suppose  $|X - \mu| \leq b$ ,  $E(X - \mu)^2 = \sigma^2$ . Then

$$E e^{\lambda(X - \mu)} \leq \exp\left(\frac{\lambda^2 \sigma^2}{2(1 - b|\lambda|)}\right) \quad \forall |\lambda| < \frac{1}{b}$$

Therefore

$$\textcircled{*} \quad P[|X - \mu| \geq t] \leq 2 e^{-\frac{t^2}{2(\sigma^2 + bt)}}$$

pf: Taylor expansion:

$$E e^{\lambda(X - \mu)} = \sum_{k=0}^{\infty} \frac{\lambda^k E(X - \mu)^k}{k!} = 1 + \frac{\lambda^2 \sigma^2}{2} + \sum_{k=3}^{\infty} \frac{\lambda^k E(X - \mu)^k}{k!}$$

$$\leq 1 + \sum_{k=2}^{\infty} \frac{\lambda^k \sigma^2}{2 \cdot 3 \cdot 4 \cdots k}$$

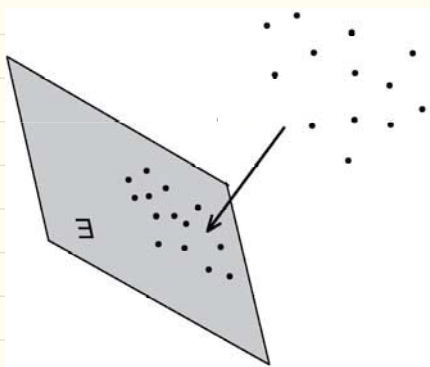
$$\leq 1 + \frac{\lambda^2 \sigma^2}{2} \cdot \frac{1}{1 - b|\lambda|} \leq \exp\left(\frac{\lambda^2 \sigma^2}{2(1 - b|\lambda|)}\right)$$

$\textcircled{*}$  Follows from Chernoff by

setting  $\lambda = \frac{t}{bt + \sigma^2} \in [0, \frac{1}{b}]$   $\square$

This is superior to Hoeffding when  $\sigma \ll b$ .

Application: Dimensionality Reduction  
 Given  $u_1, \dots, u_m \in \mathbb{R}^d$  with  $m \ll d$ ,  
 can one map  $u_1, \dots, u_m$  into a  
 lower dimensional space with low distortion?



Thm: (Johnson-Lindenstrauss)  
 Fix  $\epsilon, \delta \in (0, 1)$ , a set  $U \subseteq \mathbb{R}^d$  of  $n$   
 points and a number  $n > \frac{16 \ln(\frac{m^2}{\delta})}{\epsilon^2}$ .  
 Let  $X \in \mathbb{R}^{n \times d}$  consist of i.i.d  $N(0, 1)$  entries.  
 Then with probability  $1 - \delta$ , the map  $f(u) = \frac{1}{\sqrt{n}} Xu$   
 satisfies

$$1 - \epsilon \leq \frac{\|f(u) - f(v)\|_2^2}{\|u - v\|_2^2} \leq 1 + \epsilon \quad \forall u, v \in U.$$

pf: Observe

$$\frac{\|Xu\|_2^2}{\|u\|_2^2} = \sum_{i=1}^n \left\langle \frac{u}{\|u\|_2}, x_i \right\rangle^2$$

i.i.d.  $\mathcal{N}(0,1)$

$x_i$  is  $i$ th row of  $X$

$\Rightarrow$  Bernstein:

$$P\left[\left|\frac{\|Xu\|_2^2}{n\|u\|_2^2} - 1\right| > \varepsilon\right] \leq 2 \exp\left[-\left(\frac{n\varepsilon^2}{8} \wedge \frac{n\varepsilon}{8}\right)\right]$$

So for any  $i, j$  get  $= 2 \exp\left(-\frac{n\varepsilon^2}{8}\right) \forall 0 \leq \varepsilon \leq 1$

$$P\left[\frac{\|f(u_i - u_j)\|_2^2}{\|u_i - u_j\|_2^2} \notin [1-\varepsilon, 1+\varepsilon]\right] \leq 2e^{-n\varepsilon^2/8}$$

Take union bound over  $\binom{m}{2}$  pairs of points

$$2\binom{m}{2}e^{-n\varepsilon^2/8} \leq m^2 e^{-n\varepsilon^2/8} = \delta \quad \square$$

Question: What if  $m \rightarrow \infty$  but  $U$  has few "degrees of freedom"?



So far, we have focused on concentration of the average  $\frac{1}{n} \sum_{i=1}^n x_i$ . Often one is interested in bounds  $P[|f(x_1, \dots, x_n) - \mathbb{E} f(x_1, \dots, x_n)| > t] \leq \text{small}(n, t)$  where  $x_i$  are independent and  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is some function.

Useful insight:

As long as  $f(x_1, \dots, x_n)$  depends weakly on individual  $x_i$ , concentration holds!

Thm: (McDiarmid) Suppose that  $f: X^n \rightarrow \mathbb{R}$  has the bounded difference property:

$\exists L_1, \dots, L_n$  such that

$$|f(x_1, \dots, x_i, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq L_i \quad \forall x, x'_i \in \mathbb{R}^n$$

Then for independent  $R_i$ ,  $X = (X_1, \dots, X_n)$  have

$$P[|f(X) - \mathbb{E}f(X)| > t] \leq 2e^{-\frac{2t^2}{\|L\|_2^2}}$$

pf: We will use the Martingale method.

Define

$$y_0 = \mathbb{E}f(X) \quad \text{and} \quad y_i = \mathbb{E}[f(X) | x_1, \dots, x_i] \quad \forall i$$

Observe

$$y_i = y_0 + \sum_{j=0}^{i-1} (y_{j+1} - y_j) = y_0 + \sum_{j=0}^{i-1} D_{j+1}$$

$$\text{and } \mathbb{E}[y_i | x_1, \dots, x_{i-1}] = \mathbb{E}[f(X) | x_1, \dots, x_{i-1}] = y_{i-1}$$

$$\Rightarrow \mathbb{E}[D_{j+1} | x_1, \dots, x_j] = 0$$

$$\Rightarrow \mathbb{E}[e^{\lambda(f(X) - \mathbb{E}f(X))}] = \mathbb{E}[e^{\lambda(y_n - y_0)}]$$

$$\begin{aligned}
\mathbb{E}[e^{\lambda(f(x) - \mathbb{E}f(x))}] &= \mathbb{E}[e^{\lambda(Y_n - y_0)}] \\
&= \mathbb{E}[e^{\lambda \sum_{j=1}^n D_j}] \\
&= \mathbb{E}[e^{\lambda(Y_{n-1} - y_0)} e^{\lambda D_n}] \\
&= \mathbb{E}[e^{\lambda(Y_{n-1} - y_0)} \underbrace{\mathbb{E}[e^{\lambda D_n} | x_1, \dots, x_{n-1}]}]
\end{aligned}$$

Let  $x'$  differ from  $x$  in  $x_i$ . Then

$$\mathbb{E}[e^{\lambda D_i} | x_1, \dots, x_{i-1}] = \mathbb{E}[e^{\lambda(y_i - y_{i-1})} | x_1, \dots, x_{i-1}]$$

$$= \mathbb{E}[e^{\lambda(f(x) - f(x'))} | x_1, \dots, x_{i-1}]$$

$$\stackrel{\text{Jensen}}{\leq} \mathbb{E}[e^{\lambda(f(x) - f(x'))} | x_1, \dots, x_{i-1}]$$

$$\Rightarrow \mathbb{E}[e^{\lambda(f(x) - f(x'))} | x_1, \dots, x_{i-1}] \leq e^{\frac{\lambda^2 L_i^2}{8}}$$

↑  
bounded by  $L_i$

$$\text{So } \mathbb{E}[e^{\lambda(f(x) - \mathbb{E}f(x))}] \leq e^{\frac{\lambda^2 \|L\|^2}{8}} \quad \text{Apply Chernoff}$$

□

Recall if  $X_1, \dots, X_n$  are independent  $\sigma$ -sub Gaussian with  $\mathbb{E}X_i = \mu$ , the Hoeffding implies that  $\hat{x} = \frac{1}{n} \sum_{i=1}^n X_i$  satisfies

$$\mathbb{P}\left[|\hat{x} - \mu| \leq t\right] \geq 1 - 2\exp\left(-\frac{nt^2}{2\sigma^2}\right)$$

or equivalently

$$\mathbb{P}\left[|\hat{x} - \mu| \leq \sqrt{\frac{2\sigma^2 \ln\left(\frac{2}{\rho}\right)}{n}}\right] \geq 1 - \rho.$$

Can one achieve a similar guarantee without sub Gaussian assumption with a different estimator  $\hat{x}$ ?

Answer: yes, almost!

Thm: (Median of means)

Consider  $X \in \mathbb{R}$  with  $\mathbb{E}X = \mu$  and  $\text{Var}(X) = \sigma^2$ .  
 Let  $X_1, \dots, X_n$  be i.i.d realizations of  $X$ .  
 Subdivide into  $k = 18 \log(\frac{1}{\rho})$  bins and form the empirical means  $\hat{x}_j$  for  $j=1, \dots, k$ .

Then  $\hat{x} = \text{median}(\hat{x}_1, \dots, \hat{x}_k)$  satisfies

$$\mathbb{P}\left[|\hat{x} - \mu| \leq \sqrt{\frac{27\sigma^2 \log(\frac{1}{\rho})}{n}}\right] \geq 1 - \rho$$

pt: By Chebychev

$$\mathbb{P}[|X_i - \mu| \geq \sqrt{\frac{3\sigma^2 k}{2n}}] \leq \frac{\sigma^2}{\frac{n/k}{3\sigma^2 k / 2n}} = \frac{2}{3} \quad \forall i.$$

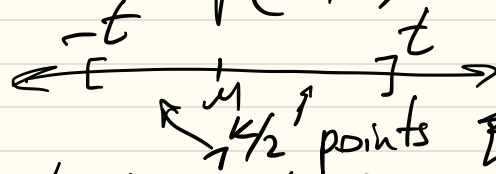
Let  $\mathbb{1}_i$  be indicator of this event

Then by Hoeffding,

$$\mathbb{P}\left[\frac{1}{k} \sum_{i=1}^k \mathbb{1}_i > \frac{1}{2}\right] \leq \exp\left(-\frac{k}{18}\right)$$

in this event,

$$\Rightarrow |\hat{x} - \mu| \leq \sqrt{\frac{3\sigma^2 k}{2n}}$$



Notice that in contrast to sub-Gaussian case,  $\hat{x}$  depends on confidence level  $\rho$ .