

# Deck 5: Some uses of scalar concentration inequalities

Math 7870: Topics in Randomized Numerical Linear Algebra

Spring 2026

Akil Narayan

## Example 1: Randomized trace estimators

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be symmetric. Our goal is to estimate the trace of  $\mathbf{A}$ ,

$$\text{tr}(\mathbf{A}) = \sum_{j \in [n]} A_{j,j} = \sum_{j \in [n]} \lambda_j(\mathbf{A}).$$

The idea for accomplishing this fairly transparent: the trace is a sum of a symmetric set of entries of a matrix.

## Example 1: Randomized trace estimators

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be symmetric. Our goal is to estimate the trace of  $\mathbf{A}$ ,

$$\text{tr}(\mathbf{A}) = \sum_{j \in [n]} A_{j,j} = \sum_{j \in [n]} \lambda_j(\mathbf{A}).$$

The idea for accomplishing this fairly transparent: the trace is a sum of a symmetric set of entries of a matrix.

So consider a centered random vector  $\mathbf{x} \in \mathbb{R}^n$ , i.e.,  $\mathbb{E}\mathbf{x} = \mathbf{0}$ . Then,

$$\mathbb{E}\mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_{i,j \in [n]} A_{i,j} \mathbb{E}(x_i x_j) = \sum_{i,j \in [n]} (\text{cov}(\mathbf{x}) \odot \mathbf{A})_{i,j}.$$

Therefore, the mean of  $\mathbf{x}^T \mathbf{A} \mathbf{x}$  is the sum of  $\mathbf{A}$ -elementwise-weighted covariance of  $\mathbf{x}$ .

## Example 1: Randomized trace estimators

Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be symmetric. Our goal is to estimate the trace of  $\mathbf{A}$ ,

$$\text{tr}(\mathbf{A}) = \sum_{j \in [n]} A_{j,j} = \sum_{j \in [n]} \lambda_j(\mathbf{A}).$$

The idea for accomplishing this fairly transparent: the trace is a sum of a symmetric set of entries of a matrix.

So consider a centered random vector  $\mathbf{x} \in \mathbb{R}^n$ , i.e.,  $\mathbb{E}\mathbf{x} = \mathbf{0}$ . Then,

$$\mathbb{E}\mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_{i,j \in [n]} A_{i,j} \mathbb{E}(x_i x_j) = \sum_{i,j \in [n]} (\text{cov}(\mathbf{x}) \odot \mathbf{A})_{i,j}.$$

Therefore, the mean of  $\mathbf{x}^T \mathbf{A} \mathbf{x}$  is the sum of  $\mathbf{A}$ -elementwise-weighted covariance of  $\mathbf{x}$ .

To recover the trace, we need a special type of random vector.

A random vector  $\mathbf{x}$  is *isotropic* if  $\mathbb{E}\mathbf{x}\mathbf{x}^T = \mathbf{I}$ .

# The (Hutchinson) randomized trace estimator

Here's the ("Hutchinson") randomized trace estimation procedure:

Input:  $N$ , an isotropic distribution for  $\mathbf{x}$ , and  $\mathbf{A}$

0. Initialize  $Z_0 = 0$ ,  $j = 0$ .
1. Generate an isotropic random vector  $\mathbf{x}$ .
2. Compute  $Z_0 \leftarrow Z_0 + \frac{1}{N} \mathbf{x}^T \mathbf{A} \mathbf{x}$ .
3. Set  $j \leftarrow j + 1$ .
4. If  $j < N$ , return to step 1. Otherwise, return  $Z_0$ .

NB: We need not have  $\mathbf{A}$ , but instead just access to operator queries  $\mathbf{x} \mapsto \mathbf{A} \mathbf{x}$ .

# The (Hutchinson) randomized trace estimator

Here's the ("Hutchinson") randomized trace estimation procedure:

Input:  $N$ , an isotropic distribution for  $\mathbf{x}$ , and  $\mathbf{A}$

0. Initialize  $Z_0 = 0, j = 0$ .
1. Generate an isotropic random vector  $\mathbf{x}$ .
2. Compute  $Z_0 \leftarrow Z_0 + \frac{1}{N} \mathbf{x}^T \mathbf{A} \mathbf{x}$ .
3. Set  $j \leftarrow j + 1$ .
4. If  $j < N$ , return to step 1. Otherwise, return  $Z_0$ .

$$Z_j \sim Z$$
$$Z = \mathbf{x}^T \mathbf{A} \mathbf{x}$$

NB: We need not have  $\mathbf{A}$ , but instead just access to operator queries  $\mathbf{x} \mapsto \mathbf{A} \mathbf{x}$ .

Let  $Z$  be the random variable  $\mathbf{x}^T \mathbf{A} \mathbf{x}$ . If  $Z_j, j \geq 1$  are iid copies of  $Z$ , the Hutchinson trace estimator produces,

$$Z_0 = \frac{1}{N} \sum_{j \in [N]} Z_j, \quad \mathbb{E} Z_0 = \mathbb{E} Z = \sum_{i,j} (\mathbb{E} \mathbf{x} \mathbf{x}^T \odot \mathbf{A})_{i,j} = \sum_{i,j} (\mathbf{I} \odot \mathbf{A})_{i,j} = \text{tr}(\mathbf{A}), \quad \text{Var} Z_0 = \frac{1}{N} \text{Var} Z.$$

## Isotropic distribution examples

There are two common choices for isotropic distributions:

- $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . Then  $\text{Var}Z = 2\|\mathbf{A}\|_F^2$ .
- $\mathbf{x}$  has iid Rademacher entries:  $x_j \sim \text{Unif}(\{-1, 1\})$ . Then  $\text{Var}Z = 4 \sum_{i < j} |A_{i,j}|^2$ .

## Isotropic distribution examples

There are two common choices for isotropic distributions:

- $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . Then  $\text{Var}Z = 2\|\mathbf{A}\|_F^2$ .
- $\mathbf{x}$  has iid Rademacher entries:  $x_j \sim \text{Unif}(\{-1, 1\})$ . Then  $\text{Var}Z = 4 \sum_{i < j} |A_{i,j}|^2$ .

To see the Rademacher result: Suppose  $\mathbf{x}$  has iid Rademacher entries. Then:

$$\begin{aligned}\text{Var}Z &= \mathbb{E} \left[ \left( \sum_{i,j} (x_i x_j) A_{i,j} \right)^2 \right] - \mathbb{E}Z^2 \\ &= \mathbb{E} \left[ \sum_{i,j,k,l} x_i x_j x_k x_l A_{i,j} A_{k,l} \right] - \left( \sum_j A_{j,j} \right)^2 \\ &= \mathbb{E} \left[ \sum_i x_i^4 A_{i,i}^2 + \sum_{i \neq k} x_i^2 x_k^2 A_{i,i} A_{k,k} + \sum_{i \neq j} x_i^2 x_j^2 A_{i,j}^2 + \sum_{i \neq j} x_i^2 x_j^2 A_{i,j}^2 \right] - \sum_i A_{i,i}^2 - \sum_{i \neq j} A_{i,i} A_{j,j} \\ &= 2 \sum_{i \neq j} A_{i,j}^2.\end{aligned}$$

## Using moments

These actually provide significant intuition about when trace estimators succeed. Consider the *intrinsic dimension* of an spd matrix  $\mathbf{A}$ :

$$\text{idim}(\mathbf{A}) = \frac{\text{tr}(\mathbf{A})}{\|\mathbf{A}\|_2} = \frac{\sum_{j \in [n]} \lambda_j(\mathbf{A})}{\max_{j \in [n]} \lambda_j(\mathbf{A})} \in [1, n]$$

The intrinsic dimension is one way to define a “continuous” version of rank/dimension of a matrix.

## Using moments

These actually provide significant intuition about when trace estimators succeed. Consider the *intrinsic dimension* of an spd matrix  $\mathbf{A}$ :

$$\text{idim}(\mathbf{A}) = \frac{\text{tr}(\mathbf{A})}{\|\mathbf{A}\|_2} = \frac{\sum_{j \in [n]} \lambda_j(\mathbf{A})}{\max_{j \in [n]} \lambda_j(\mathbf{A})} \in [1, n]$$

The intrinsic dimension is one way to define a “continuous” version of rank/dimension of a matrix.

Chaining together the previous results and using  $\|\mathbf{A}\|_F^2 \leq \|\mathbf{A}\|_2 \text{tr}(\mathbf{A})$ , we conclude for either distribution above and Chebyshev’s inequality:

$$\Pr(|Z_0 - \text{tr}(\mathbf{A})| > \epsilon \text{tr}(\mathbf{A})) \leq \frac{2}{N\epsilon^2 \text{idim}(\mathbf{A})}.$$

I.e., this method works well for matrices with large intrinsic dimension.

# Concentration

Of course, we know we can do better than Chebyshev's inequality!

Given  $\mathbf{A}$ :

$$|Z - \mathbb{E}Z| \leq \sum_{i \neq j} |A_{i,j}| \quad \text{wp1.}$$

So  $Z - \mathbb{E}Z$  is a bounded random variable.

# Concentration

Of course, we know we can do better than Chebyshev's inequality!

Given  $\mathbf{A}$ :

$$|Z - \mathbb{E}Z| \leq \sum_{i \neq j} |A_{i,j}| \quad \text{wp1.}$$

So  $Z - \mathbb{E}Z$  is a bounded random variable.

Applying Hoeffding's inequality for bounded random variables: Taking  $N$  samples, we have

$$\Pr(|Z_0 - \text{tr}(\mathbf{A})| \gtrsim \epsilon \text{tr}(\mathbf{A})) \leq 2 \exp\left(-\frac{N\epsilon^2}{\text{Var}Z}\right),$$

I.e., if  $N \gtrsim r \log N \frac{\text{Var}Z}{\epsilon^2}$ , then this succeeds with probability  $\gtrsim 1 - N^{-r}$ , which is *much* better than Chebyshev's inequality!

## Why are trace estimators useful?

If I have easy access to entries of  $\mathbf{A}$ , it's easy to exactly compute  $\text{tr}(\mathbf{A})$ .

These estimators are useful when the entries of  $\mathbf{A}$  are not easily available, but instead when  $\mathbf{x} \mapsto \mathbf{Ax}$  is available.

## Why are trace estimators useful?

If I have easy access to entries of  $\mathbf{A}$ , it's easy to exactly compute  $\text{tr}(\mathbf{A})$ .

These estimators are useful when the entries of  $\mathbf{A}$  are not easily available, but instead when  $\mathbf{x} \mapsto \mathbf{A}\mathbf{x}$  is available.

For example:

- When the entries of  $\mathbf{A}$  are simply not available.
- When computing  $\text{tr} \left[ (z\mathbf{I} - \mathbf{A})^{-1} \right]$
- When computing  $\text{tr}f(\mathbf{A})$ , often using truncated Taylor approximations of  $f$  (i.e., polynomials in  $\mathbf{A}$ ).

However, in any of these situations  $\mathbf{x} \mapsto \mathbf{A}\mathbf{x}$  could be used to recover entries of  $\mathbf{A}$  by choosing  $\mathbf{x}$  appropriately.

Therefore, a randomized trace estimator would only be useful if we can choose  $N$  much smaller than the dimension of  $\mathbf{A}$ .

# Gaussian sampling

We hope an exponential concentration argument, similar to the Rademacher case, works for Gaussian isotropic vectors.

However, note that if  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ , then

$$Z = \mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_i A_{i,i} x_i^2 + \sum_{i \neq j} A_{i,j} x_{i,j}^2,$$

# Gaussian sampling

We hope an exponential concentration argument, similar to the Rademacher case, works for Gaussian isotropic vectors.

However, note that if  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ , then

$$Z = \mathbf{x}^T \mathbf{A} \mathbf{x} = \sum_i A_{i,i} x_i^2 + \sum_{i \neq j} A_{i,j} x_{i,j}^2,$$

The important observation is that, at least part of this quantity,  $x_i^2$ , is *not* a normal random variable. Indeed:

$$\Pr(x_i^2 \geq t) = \Pr(|x_i| \geq \sqrt{t}) \sim e^{-t},$$

i.e., this random variable is not sub-Gaussian.

(It's exactly a Gamma/exponential random variable.)

To handle these types of random variables, we require some extra study of concentration.

## Back to MGF's

When studying concentration in the simplest setting, we have the setup:

$$Z_n = \frac{1}{n} S_n := \frac{1}{n} \sum_{i \in [n]} X_i \quad X_i \stackrel{\text{iid}}{\sim} X, \quad \mathbb{E}X = 0$$

And we want to know tail probabilities of  $Z_n$ .

Recall, for any  $s > 0$ :

$$M_X(s) = \mathbb{E} e^{sX}$$

$$\Pr(Z_n \geq t) = \Pr(S_n \geq nt) \leq e^{-nst} [M_X^n(s)]$$

While we've introduced this for  $s > 0$ , taking  $s < 0$  bounds  $\Pr(Z_n \leq t)$ , so that considering all  $s \in \mathbb{R}$  bounds both tails.

$$\Pr(Z_n \leq t) \xrightarrow{s < 0} \Pr(sZ_n \geq st)$$

## Back to MGF's

When studying concentration in the simplest setting, we have the setup:

$$Z_n = \frac{1}{n} S_n := \frac{1}{n} \sum_{i \in [n]} X_i \quad X_i \stackrel{\text{iid}}{\sim} X, \quad \mathbb{E}X = 0$$

And we want to know tail probabilities of  $Z_n$ .

Recall, for any  $s > 0$ :

$$\Pr(Z_n \geq t) = \Pr(S_n \geq nt) \leq e^{-nst} [M_X^n(s)]$$

While we've introduced this for  $s > 0$ , taking  $s < 0$  bounds  $\Pr(Z_n \leq t)$ , so that considering all  $s \in \mathbb{R}$  bounds both tails.

In particular, this implies:

$$\begin{aligned} \Pr(Z_n \geq t) &\leq \inf_{s \in \mathbb{R}} e^{-nst} M_X^n(s) \\ &= \exp \left( n \left( \inf_s \left( -st + \log M_X(s) \right) \right) \right) \\ &= \exp \left( -n \sup_s \left[ st - \log M_X(s) \right] \right). \end{aligned}$$

*NB:  $M_X$  need not be defined on all  $\mathbb{R}$ !*

## CGF's and rate functions

With a little rearrangement, we've shown the following:

$$\frac{1}{n} \log \Pr(Z_n \geq t) \leq -I(t), \quad I(t) = \sup_s (st - \log M_X(s)).$$

In the language of large deviations,  $I(t)$  is the *rate function*, which characterizes tail probabilities.

## CGF's and rate functions

With a little rearrangement, we've shown the following:

$$\frac{1}{n} \log \Pr(Z_n \geq t) \leq -I(t), \quad I(t) = \sup_s (st - \log M_X(s)).$$

In the language of large deviations,  $I(t)$  is the *rate function*, which characterizes tail probabilities.

Equivalently, define the function  $K_X$ ,

$$K_X(s) := \log M_X(s) = \log \mathbb{E} e^{sX},$$

which is the *cumulant generating function* (CGF) of  $X$ .

# CGF's and rate functions

With a little rearrangement, we've shown the following:

$$\frac{1}{n} \log \Pr(Z_n \geq t) \leq -I(t), \quad I(t) = \sup_s (st - \log M_X(s)).$$

In the language of large deviations,  $I(t)$  is the *rate function*, which characterizes tail probabilities.

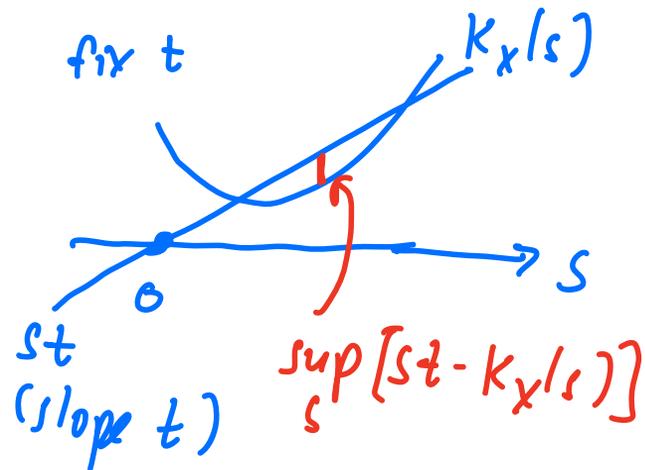
Equivalently, define the function  $K_X$ ,

$$K_X(s) := \log M_X(s) = \log \mathbb{E}e^{sX},$$

which is the *cumulant generating function* (CGF) of  $X$ .

Then  $I(t)$  is related to  $K_X(s)$  by:

$$I_\bullet(t) = \sup_s (st - K_X(s))$$



# The Legendre-Fenchel transform

$$\log \Pr(Z_n \geq t) \leq -I(t),$$

$$I(t) = \sup_s (st - K_X(s))$$

The operation connecting  $I$  to  $K$  itself has a name:

Let  $f : \mathbb{R} \rightarrow \mathbb{R}$ . The function  $f^*$  defined as,

$$f^*(t) := \sup_s (st - f(s)),$$

is called the *Legendre-Fenchel transform/dual* or *convex conjugate* of  $f$ .

# The Legendre-Fenchel transform

$$\frac{1}{n} \log \Pr(Z_n \geq t) \leq -I(t),$$

$$I(t) = \sup_s (st - K_X(s))$$

The operation connecting  $I$  to  $K$  itself has a name:

Let  $f : \mathbb{R} \rightarrow \mathbb{R}$ . The function  $f^*$  defined as,

$$f^*(t) := \sup_s (st - f(s)),$$

is called the *Legendre-Fenchel transform/dual* or *convex conjugate* of  $f$ .

So with all this new language:

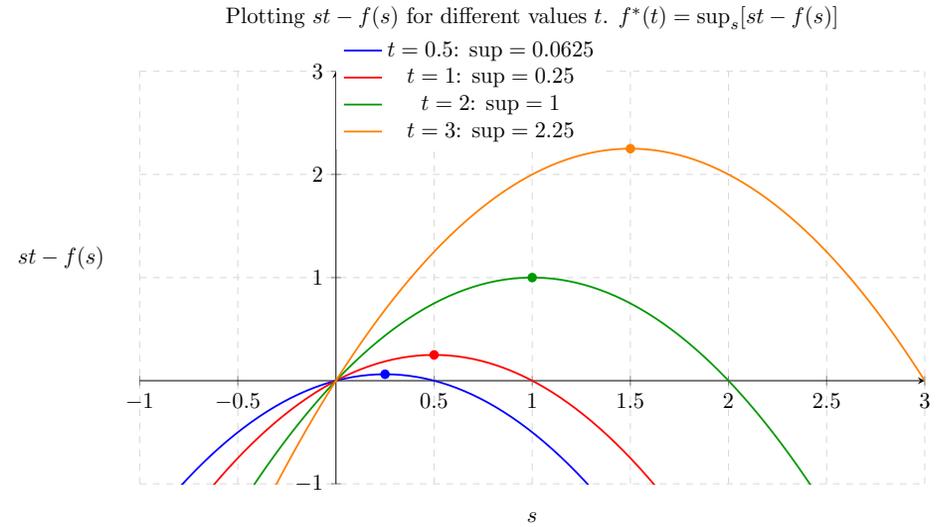
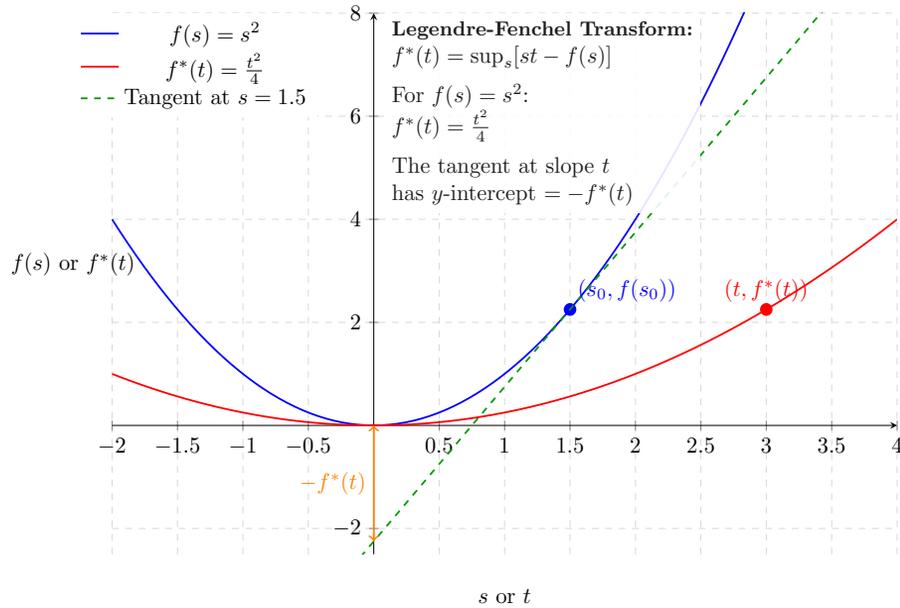
The rate function (corresponding to the tail probability we seek to compute) is the convex conjugate of the CGF of  $X$ :  $I(t) = K_X^*(t)$ .

Or:

The tail probability is the exponential of the negative convex conjugate of the CGF of  $X$ :

$$\Pr(Z_n \geq t) \leq \exp(-nI(t))$$

# Legendre-Fenchel visualization



$$f(s) = as^2, \quad a > 0$$

$$f^*(t) = \sup_s (st - f(s)) = \sup_s (st - as^2)$$

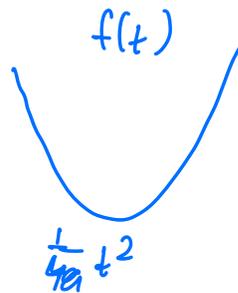
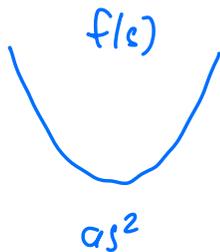
$$g(s) = st - as^2$$

$$g'(s) = t - 2as$$

$$\Rightarrow g'(s) = 0 \Rightarrow s = \frac{t}{2a}$$

$$f^*(t) = \sup_s g(s) = \sup_s \left[ \frac{t}{2a} \cdot t - a \cdot \left(\frac{t}{2a}\right)^2 \right]$$

$$= \frac{t^2}{2a} - \frac{t^2}{4a} = \frac{t^2}{4a}$$



If  $f$  is smooth: fix  $t$ , solve  $\sup_s (st - f(s))$

$$s^* : f'(s^*) = t$$

$$s^* = (f')^{-1}(t)$$

$$f^*(t) = s^*t - f(s^*)$$

# Legendre-Fenchel transform

$$f^*(t) := \sup_s (st - f(s)), \quad I(t) = K_X^*(t), \quad \Pr(Z_n \geq t) \leq \exp(-nI(t))$$

The body of knowledge on Legendre-Fenchel transforms is large. Here are some punch lines that are somewhat relevant for us:

- Convex conjugates of convex functions are convex.
- CGF's are convex functions  $\implies$  Rate functions are convex.
- The Legendre-Fenchel transform is an involution on convex functions:  $f^{**} = f$ .

# Legendre-Fenchel transform

$$f^*(t) := \sup_s (st - f(s)), \quad I(t) = K_X^*(t), \quad \Pr(Z_n \geq t) \leq \exp(-nI(t))$$

The body of knowledge on Legendre-Fenchel transforms is large. Here are some punch lines that are somewhat relevant for us:

- Convex conjugates of convex functions are convex.
- CGF's are convex functions  $\implies$  Rate functions are convex.
- The Legendre-Fenchel transform is an involution on convex functions:  $f^{**} = f$ .
- We know the convex conjugate for the MGF of many named ("classical") distributions.
- $X \sim \mathcal{N}(0, \sigma^2) \implies K_X(s) = \frac{1}{2}\sigma^2 s^2 \implies K_X^*(t) = \frac{1}{2\sigma^2} t^2$ .

# Legendre-Fenchel transform

$$f^*(t) := \sup_s (st - f(s)), \quad I(t) = K_X^*(t), \quad \Pr(Z_n \geq t) \leq \exp(-nI(t))$$

The body of knowledge on Legendre-Fenchel transforms is large. Here are some punch lines that are somewhat relevant for us:

- Convex conjugates of convex functions are convex.
- CGF's are convex functions  $\implies$  Rate functions are convex.
- The Legendre-Fenchel transform is an involution on convex functions:  $f^{**} = f$ .
- We know the convex conjugate for the MGF of many named ("classical") distributions.
- $X \sim \mathcal{N}(0, \sigma^2) \implies K_X(s) = \frac{1}{2}\sigma^2 s^2 \implies K_X^*(t) = \frac{1}{2\sigma^2} t^2$ .
- So: If  $X$  is  $\mathcal{N}(0, \sigma^2)$ , then:

$$I(t) = -\frac{1}{n} \log \Pr(Z_n \geq t) \geq K_X^*(t) = \frac{t^2}{2\sigma^2} \implies \Pr(Z_n \geq t) \leq \exp\left(-\frac{nt^2}{2\sigma^2}\right).$$

## Exponential distributions

We can apply this general technique to a new, relevant situation: If  $X \sim \mathcal{N}(0, 1)$ , then  $X^2$  has tail probabilities behaving like  $f_{X^2}(x) \sim \exp(-|x|)$ .

So, suppose  $Y$  has density  $f_Y(y) = \frac{1}{2} \exp(-|y|)$ . Let's compute the tail probability for an iid sum:

$$Z_n = \frac{1}{n} \sum_{i \in [n]} Y_i, \quad Y_i \stackrel{\text{iid}}{\sim} Y.$$

## Exponential distributions

We can apply this general technique to a new, relevant situation: If  $X \sim \mathcal{N}(0, 1)$ , then  $X^2$  has tail probabilities behaving like  $f_{X^2}(x) \sim \exp(-|x|)$ .

So, suppose  $Y$  has density  $f_Y(y) = \frac{1}{2} \exp(-|y|)$ . Let's compute the tail probability for an iid sum:

$$Z_n = \frac{1}{n} \sum_{i \in [n]} Y_i, \quad Y_i \stackrel{\text{iid}}{\sim} Y.$$

As before, we have:

$$I(t) := K_Y^*(t), \quad \Pr(Z_n \geq t) \leq \exp(-nI(t)).$$

$$\begin{aligned} M_Y(s) &= \frac{1}{2} \int_{\mathbb{R}} e^{sy} e^{-|y|} dy = \frac{1}{2} \int_{-\infty}^0 e^{(s+1)y} dy + \frac{1}{2} \int_0^{\infty} e^{(s-1)y} dy \\ &= \frac{1}{2} \left[ \frac{1}{s+1} - \frac{1}{s-1} \right] = \frac{1}{2} \left( \frac{s-1-s-1}{s^2-1} \right) = \frac{1}{1-s^2} \end{aligned}$$

## Exponential distributions

We can apply this general technique to a new, relevant situation: If  $X \sim \mathcal{N}(0, 1)$ , then  $X^2$  has tail probabilities behaving like  $f_{X^2}(x) \sim \exp(-|x|)$ .

So, suppose  $Y$  has density  $f_Y(y) = \frac{1}{2} \exp(-|y|)$ . Let's compute the tail probability for an iid sum:

$$Z_n = \frac{1}{n} \sum_{i \in [n]} Y_i, \quad Y_i \stackrel{\text{iid}}{\sim} Y.$$

As before, we have:

$$I(t) := K_Y^*(t), \quad \Pr(Z_n \geq t) \leq \exp(-nI(t)).$$

We directly compute that  $M_Y(s) = \frac{1}{1-s^2}$ ,  $|s| < 1$ .

Of particular note here:  $M_Y(s)$  is defined only on a subset of  $\mathbb{R}$ . It does not exist for  $|s| \geq 1$ .

$$\Rightarrow K_Y(s) = -\log(1-s^2)$$

$$K_Y^*(t) = \sup_s (st - K_Y(s)), \text{ define } g_t(s) = st + \log(1-s^2)$$

$$\text{Compute CP of } g_t(s): g_t'(s) = 1 - \frac{2s}{1-s^2} = 0$$

$$1 - s^2 = \frac{2s}{t}$$

$$1 = s^2 + \frac{2s}{t} + \frac{1}{t^2} - \frac{1}{t^2}$$

$$1 + \frac{1}{t^2} = \left(s + \frac{1}{t}\right)^2$$

$$s = -\frac{1}{t} \pm \sqrt{1 + \frac{1}{t^2}}$$

$$|s| < 1 \Rightarrow s^* = -\frac{1}{t} + \sqrt{1 + \frac{1}{t^2}}$$

$$K_y^*(t) = s^* t - K_y(s^*) \\ = s^* t + \log(1 - (s^*)^2), \quad s^* = -\frac{1}{t} + \sqrt{1 + \frac{1}{t^2}}$$

## Exponential distributions

We can apply this general technique to a new, relevant situation: If  $X \sim \mathcal{N}(0, 1)$ , then  $X^2$  has tail probabilities behaving like  $f_{X^2}(x) \sim \exp(-|x|)$ .

So, suppose  $Y$  has density  $f_Y(y) = \frac{1}{2} \exp(-|y|)$ . Let's compute the tail probability for an iid sum:

$$Z_n = \frac{1}{n} \sum_{i \in [n]} Y_i, \quad Y_i \stackrel{\text{iid}}{\sim} Y.$$

As before, we have:

$$I(t) := K_Y^*(t), \quad \Pr(Z_n \geq t) \leq \exp(-nI(t)).$$

We directly compute that  $M_Y(s) = \frac{1}{1-s^2}$ ,  $|s| < 1$ .

Of particular note here:  $M_Y(s)$  is defined only on a subset of  $\mathbb{R}$ . It does not exist for  $|s| \geq 1$ .

We therefore have  $K_Y(s) = -\log(1-s^2)$ , which is smooth (and convex), and using

$$K_Y^*(t) = ts^* - K_Y(s^*), \quad s^*(t) = (K'_Y)^{-1}(t),$$

we identify

$$K_Y^*(t) = \log\left(\frac{2}{e}\right) + \sqrt{1+t^2} + \log\left(\frac{1}{t^2} \left(\sqrt{1+t^2} - 1\right)\right).$$

## Exponential tail probabilities

We can now identify tail probability behavior:

$$Z_n = \frac{1}{n} \sum_{i \in [n]} Y_i, \quad Y_i \stackrel{\text{iid}}{\sim} Y, \quad f_Y(y) = \frac{1}{2} \exp(-|y|)$$

$$K_Y^*(t) = \log\left(\frac{2}{e}\right) + \sqrt{1+t^2} + \log\left(\frac{1}{t^2} \left(\sqrt{1+t^2} - 1\right)\right).$$

## Exponential tail probabilities

We can now identify tail probability behavior:

$$Z_n = \frac{1}{n} \sum_{i \in [n]} Y_i, \quad Y_i \stackrel{\text{iid}}{\sim} Y, \quad f_Y(y) = \frac{1}{2} \exp(-|y|)$$

$$K_Y^*(t) = \log\left(\frac{2}{e}\right) + \sqrt{1+t^2} + \log\left(\frac{1}{t^2} \left(\sqrt{1+t^2} - 1\right)\right).$$

Through Taylor series ( $\sqrt{1+x^2} \sim 1 + x^2/2$  for small  $x$ ), one can show:

$$t \ll 1 \implies K_Y^*(t) \sim -1 + \frac{1}{2}t^2$$

$$t \gg 1 \implies K_Y^*(t) \sim t.$$

# Exponential tail probabilities

We can now identify tail probability behavior:

$$Z_n = \frac{1}{n} \sum_{i \in [n]} Y_i, \quad Y_i \stackrel{\text{iid}}{\sim} Y, \quad f_Y(y) = \frac{1}{2} \exp(-|y|)$$

$$K_Y^*(t) = \log\left(\frac{2}{e}\right) + \sqrt{1+t^2} + \log\left(\frac{1}{t^2} \left(\sqrt{1+t^2} - 1\right)\right).$$

Through Taylor series ( $\sqrt{1+x^2} \sim 1 + x^2/2$  for small  $x$ ), one can show:

$$t \ll 1 \implies K_Y^*(t) \sim -1 + \frac{1}{2}t^2$$

$$t \gg 1 \implies K_Y^*(t) \sim t.$$

i.e.:

$$t \ll 1 \implies \Pr(Z_n \geq t) \sim \exp\left(-\frac{nt^2}{2}\right)$$

$$t \gg 1 \implies \Pr(Z_n \geq t) \sim \exp(-nt).$$

And one can show that these bounds can be stitched together to cover all  $t > 0$ .

## Sub-exponential distributions

In order to state the generalization of this result, we consider the appropriate class of exponential-like distributions.

**Definition 1** (Sub-exponential distributions). *A random variable  $X$  has sub-exponential distribution if there is a constant  $C > 0$  such that:*

$$\Pr(|X| \geq t) \leq 2 \exp(-Ct).$$

## Sub-exponential distributions

In order to state the generalization of this result, we consider the appropriate class of exponential-like distributions.

**Definition 1** (Sub-exponential distributions). *A random variable  $X$  has sub-exponential distribution if there is a constant  $C > 0$  such that:*

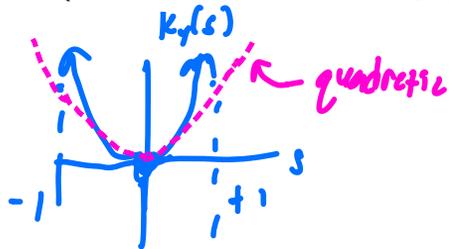
$$\Pr(|X| \geq t) \leq 2 \exp(-Ct).$$

As expected, there is an equivalent definition that is more useful for tail probabilities:

**Theorem** (Equivalent definitions of sub-exponential distributions). *The definition above for a sub-exponential distribution  $X$  is equivalent to: there exists a  $c > 0$  such that*

$$K_X(s) \leq c^2 s^2, \quad |s| \leq \frac{1}{c}$$

(Recall for  $f_X \sim \exp(-|x|)$ , we had  $K_X(s) \leq s^2$  for  $|s| < 1$ .)



"

$$- \log(1-s^2)$$

$|s| < c, c < 1$

## Sub-exponential distributions

In order to state the generalization of this result, we consider the appropriate class of exponential-like distributions.

**Definition 1** (Sub-exponential distributions). *A random variable  $X$  has sub-exponential distribution if there is a constant  $C > 0$  such that:*

$$\Pr(|X| \geq t) \leq 2 \exp(-Ct).$$

As expected, there is an equivalent definition that is more useful for tail probabilities:

**Theorem** (Equivalent definitions of sub-exponential distributions). *The definition above for a sub-exponential distribution  $X$  is equivalent to: there exists a  $c > 0$  such that*

$$K_X(s) \leq c^2 s^2, \quad |s| \leq \frac{1}{c}$$

(Recall for  $f_X \sim \exp(-|x|)$ , we had  $K_X(s) \leq s^2$  for  $|s| < 1$ .)

Sub-Gaussian distributions are sub-exponential distributions.

In particular: Sub-Gaussian distributions have  $K_X(s) \sim s^2$  for all  $s \in \mathbb{R}$ . Sub-exponential distributions require this only in a neighborhood of the origin.

# Bernstein's inequality

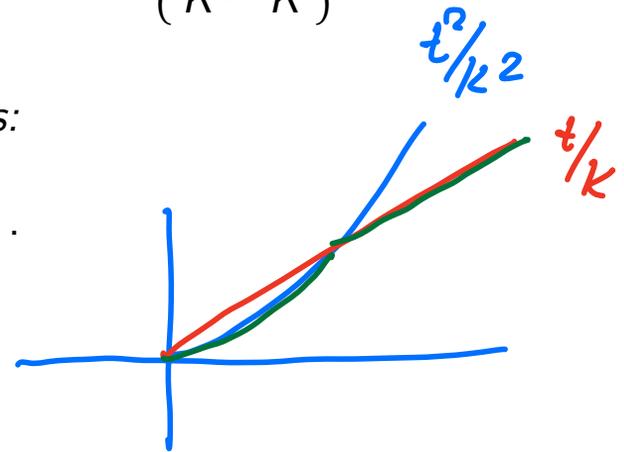
Sub-exponential distributions allows us to codify our previous observation that such random variables have tail probabilities behaving like  $e^{-t^2}$  when  $t$  is small, but  $e^{-t}$  when  $t$  is large.

**Theorem** (Bernstein's inequality). *Let  $X_i, i \in [n]$  be independent and centered sub-exponential random variables. Then there is a  $c > 0$  such that:*

$$\Pr \left( \left| \frac{1}{n} \sum_{i \in [n]} X_i \right| \geq t \right) \leq 2 \exp(-c n h_X(t)), \quad h_X(t) = \min \left\{ \frac{t^2}{K^2}, \frac{t}{K} \right\},$$

where  $K$  is the maximum sub-exponential norm of the  $X_i$ , defined as:

$$K = \max_{i \in [n]} \inf \left\{ t > 0 \mid \mathbb{E} \exp(|X_i|/t) \leq 2 \right\}.$$



## Bernstein's inequality

Sub-exponential distributions allows us to codify our previous observation that such random variables have tail probabilities behaving like  $e^{-t^2}$  when  $t$  is small, but  $e^{-t}$  when  $t$  is large.

**Theorem** (Bernstein's inequality). *Let  $X_i, i \in [n]$  be independent and centered sub-exponential random variables. Then there is a  $c > 0$  such that:*

$$\Pr \left( \left| \frac{1}{n} \sum_{i \in [n]} X_i \right| \geq t \right) \leq 2 \exp(-c n h_X(t)), \quad h_X(t) = \min \left\{ \frac{t^2}{K^2}, \frac{t}{K} \right\},$$

where  $K$  is the maximum sub-exponential norm of the  $X_i$ , defined as:

$$K = \max_{i \in [n]} \inf \{ t > 0 \mid \mathbb{E} \exp(|X|/t) \leq 2 \}.$$

Bernstein's inequality says that for sub-exponential random variables:

- Deviations close to the mean ("small deviations") have a Gaussian, CLT-type probability
- Deviations far from the mean ("large deviations") have an exponential-type probability.

## Another Bernstein inequality

There are a few versions of Bernstein's inequality. Here's one that applies to bounded random variables:

**Theorem** (Bernstein's inequality, II). *Let  $X_i$ ,  $i \in [n]$  be independent, centered, and bounded random variables, with  $|X_i| \leq K$  for all  $i \in [n]$ . Then:*

$$\Pr \left( \left| \sum_{i \in [n]} X_i \right| \geq t \right) \leq 2 \exp \left( - \frac{t^2/2}{\sum_{i \in [n]} \mathbb{E} X_i^2 + K t/2} \right).$$

This form explicitly shows for bounded random variables the tail probabilities are  $\sim \exp(-t)$  for large  $t$  and  $\sim \exp(-t^2)$  for small  $t$ .

## Another Bernstein inequality

There are a few versions of Bernstein's inequality. Here's one that applies to bounded random variables:

**Theorem** (Bernstein's inequality, II). *Let  $X_i, i \in [n]$  be independent, centered, and bounded random variables, with  $|X_i| \leq K$  for all  $i \in [n]$ . Then:*

$$\Pr \left( \left| \sum_{i \in [n]} X_i \right| \geq t \right) \leq 2 \exp \left( - \frac{t^2/2}{\sum_{i \in [n]} \mathbb{E} X_i^2 + K t/2} \right).$$

This form explicitly shows for bounded random variables the tail probabilities are  $\sim \exp(-t)$  for large  $t$  and  $\sim \exp(-t^2)$  for small  $t$ .

Without belaboring the point: now that we know that “Gaussian<sup>2</sup>” random variables have exponential concentration, we could use this to get a exponential success probability on the isotropic Gaussian randomized trace estimator.

## Example 2: The Johnson-Lindenstrauss Lemma

A very nice, “straightforward” corollary of all this is a celebrated result in dimension reduction:

**Lemma** (Johnson-Lindenstrauss, Gaussian version). *Let  $\mathbf{X} \in \mathbb{R}^{n \times m}$  be a deterministic matrix,*

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m].$$

*Let  $\mathbf{A} \in \mathbb{R}^{k \times n}$  have iid standard normal entries. Define:*

$$\mathbf{Y} = \frac{1}{\sqrt{k}} \mathbf{A} \mathbf{X} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_m] \in \mathbb{R}^{k \times m}$$

*Then:*

$$k \gtrsim \frac{(r+2) \log m}{\epsilon^2} \quad \text{w/ probability } \gtrsim 1 - m^{-r} \quad \implies \quad (1 - \epsilon) \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \leq \|\mathbf{y}_i - \mathbf{y}_j\|_2^2 \leq (1 + \epsilon) \|\mathbf{x}_i - \mathbf{x}_j\|_2^2,$$

*for all  $i, j \in [m]$ .*

I.e., with dimension complexity *logarithmic* in the number of samples, arbitrary point clouds can be near-isometrically projected.

# Johnson-Lindenstrauss proof sketch, I

We have all the necessary ingredients to prove this directly:

First, let  $\mathbf{z} \in \mathbb{R}^n$ ,  $\|\mathbf{z}\|_2 = 1$ , be deterministic and arbitrary. We'll first show that  $\mathbf{y} := \frac{1}{\sqrt{k}}\mathbf{A}\mathbf{z}$  has norm similar to  $\mathbf{z}$ .

# Johnson-Lindenstrauss proof sketch, I

We have all the necessary ingredients to prove this directly:

First, let  $\mathbf{z} \in \mathbb{R}^n$ ,  $\|\mathbf{z}\|_2 = 1$ , be deterministic and arbitrary. We'll first show that  $\mathbf{y} := \frac{1}{\sqrt{k}}\mathbf{A}\mathbf{z}$  has norm similar to  $\mathbf{z}$ .

We have that  $\mathbf{y}$  is a centered, normally distributed random vector, with,

$$\begin{aligned}\text{Cov}(\mathbf{y})_{i,j} &= \mathbb{E}y_i y_j \\ &= \frac{1}{k} \sum_{q,\ell \in [n]} \mathbb{E}A_{i,q}A_{j,\ell}z_q z_\ell \\ &= \frac{1}{k} \delta_{i,j} \sum_{q \in [n]} z_q^2 = \frac{1}{k} \delta_{i,j}.\end{aligned}$$

# Johnson-Lindenstrauss proof sketch, I

We have all the necessary ingredients to prove this directly:

First, let  $\mathbf{z} \in \mathbb{R}^n$ ,  $\|\mathbf{z}\|_2 = 1$ , be deterministic and arbitrary. We'll first show that  $\mathbf{y} := \frac{1}{\sqrt{k}}\mathbf{A}\mathbf{z}$  has norm similar to  $\mathbf{z}$ .

We have that  $\mathbf{y}$  is a centered, normally distributed random vector, with,

$$\begin{aligned}\text{Cov}(\mathbf{y})_{i,j} &= \mathbb{E}y_i y_j \\ &= \frac{1}{k} \sum_{q,\ell \in [n]} \mathbb{E}A_{i,q}A_{j,\ell}z_q z_\ell \\ &= \frac{1}{k} \delta_{i,j} \sum_{q \in [n]} z_q^2 = \frac{1}{k} \delta_{i,j}.\end{aligned}$$

Therefore,  $\|\mathbf{y}\|_2^2 = \sum_{j \in [k]} y_j^2$  is a sum of  $k$   $\mathcal{N}(0, 1/k)$ -squared iid random variables, each of which is sub-exponential.

## Johnson-Lindenstrauss proof sketch, II

Let  $W_i = y_i^2$ , so that  $\mathbb{E} \|\mathbf{y}_i\|_2^2 = \mathbb{E} \sum_{i \in [k]} W_i = 1$ , with  $W_i$  iid subexponential, having subexponential norm  $\sim 1/\sqrt{k}$ .

## Johnson-Lindenstrauss proof sketch, II

Let  $W_i = y_i^2$ , so that  $\mathbb{E}\|\mathbf{y}_i^2\|_2 = \mathbb{E}\sum_{i \in [k]} W_i = 1$ , with  $W_i$  iid subexponential, having subexponential norm  $\sim 1/\sqrt{k}$ .

We are interested in small deviations from the mean. For small  $\epsilon > 0$ , we have our concentration estimate:

$$\Pr\left(\left|\sum_{i \in [k]} W_i - 1\right| \geq \epsilon\right) \lesssim \exp(-k\epsilon^2)$$

$$\exp\left(-\frac{\epsilon^2}{(1/\sqrt{k})^2}\right)$$

$$"t = \epsilon"$$

$$"K = 1/\sqrt{k}"$$

## Johnson-Lindenstrauss proof sketch, II

Let  $W_i = y_i^2$ , so that  $\mathbb{E}\|\mathbf{y}_i^2\|_2 = \mathbb{E}\sum_{i \in [k]} W_i = 1$ , with  $W_i$  iid subexponential, having subexponential norm  $\sim 1/\sqrt{k}$ .

We are interested in small deviations from the mean. For small  $\epsilon > 0$ , we have our concentration estimate:

$$\Pr\left(\left|\sum_{i \in [k]} W_i - 1\right| \geq \epsilon\right) \lesssim \exp(-k\epsilon^2)$$

I.e., the following occurs with probability at least  $1 - ce^{-k\epsilon^2}$ :

$$\uparrow \cdot (1 - \epsilon) \leq \|\mathbf{y}\|_2^2 \leq (1 + \epsilon) \cdot \downarrow$$

with  $\|\mathbf{z}\|_2^2 = 1$ .

## Johnson-Lindenstrauss proof sketch, III

We've shown that with probability  $1 - ce^{-k\epsilon^2}$ :

$$(1 - \epsilon) \leq \left\| \frac{\frac{1}{\sqrt{k}} \mathbf{A} \mathbf{z}}{\|\mathbf{z}\|_2} \right\|_2^2 \leq (1 + \epsilon),$$

Now choose  $\mathbf{z} = \mathbf{x}_i - \mathbf{x}_j$  for all  $i, j \in [m]$  with  $i \neq j$ . Then:

$$(1 - \epsilon) \leq \left\| \frac{\frac{1}{\sqrt{k}} \mathbf{A}(\mathbf{x}_i - \mathbf{x}_j)}{\|\mathbf{x}_i - \mathbf{x}_j\|_2} \right\|_2^2 \leq (1 + \epsilon),$$

holds, by the union bound, for all  $i \neq j$  with probability at least  $1 - cm^2 e^{-k\epsilon^2}$ .

## Johnson-Lindenstrauss proof sketch, III

We've shown that with probability  $1 - ce^{-k\epsilon^2}$ :

$$(1 - \epsilon) \leq \left\| \frac{\frac{1}{\sqrt{k}} \mathbf{A} \mathbf{z}}{\|\mathbf{z}\|_2} \right\|_2^2 \leq (1 + \epsilon),$$

Now choose  $\mathbf{z} = \mathbf{x}_i - \mathbf{x}_j$  for all  $i, j \in [m]$  with  $i \neq j$ . Then:

$$(1 - \epsilon) \leq \left\| \frac{\frac{1}{\sqrt{k}} \mathbf{A}(\mathbf{x}_i - \mathbf{x}_j)}{\|\mathbf{x}_i - \mathbf{x}_j\|_2} \right\|_2^2 \leq (1 + \epsilon),$$

holds, by the union bound, for all  $i \neq j$  with probability at least  $1 - cm^2 e^{-k\epsilon^2}$ .

Pick  $k \gtrsim (r + 2) \frac{\log m}{\epsilon^2}$  to ensure this happens with probability at least,

$$1 - cm^2 \exp(-(r + 2) \log m) = 1 - cm^{-r}.$$

## Johnson-Lindenstrauss proof sketch, III

We've shown that with probability  $1 - ce^{-k\epsilon^2}$ :

$$(1 - \epsilon) \leq \left\| \frac{\frac{1}{\sqrt{k}} \mathbf{A} \mathbf{z}}{\|\mathbf{z}\|_2} \right\|_2^2 \leq (1 + \epsilon),$$

Now choose  $\mathbf{z} = \mathbf{x}_i - \mathbf{x}_j$  for all  $i, j \in [m]$  with  $i \neq j$ . Then:

$$(1 - \epsilon) \leq \left\| \frac{\frac{1}{\sqrt{k}} \mathbf{A}(\mathbf{x}_i - \mathbf{x}_j)}{\|\mathbf{x}_i - \mathbf{x}_j\|_2} \right\|_2^2 \leq (1 + \epsilon),$$

holds, by the union bound, for all  $i \neq j$  with probability at least  $1 - cm^2 e^{-k\epsilon^2}$ .

Pick  $k \gtrsim (r + 2) \frac{\log m}{\epsilon^2}$  to ensure this happens with probability at least,

$$1 - cm^2 \exp(-(r + 2) \log m) = 1 - cm^{-r}.$$

NB: The whole thing rested on us being able to bound tail probabilities of sub-exponential random variables.