

# Sketching Algorithms

Jelani Nelson

December 3, 2020



# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Probability Review . . . . .	5
<b>2</b>	<b>Counting Problems</b>	<b>13</b>
2.1	Approximate counting . . . . .	13
2.1.1	Analysis of Morris' algorithm . . . . .	13
2.1.2	Morris+ . . . . .	14
2.1.3	Morris++ . . . . .	15
2.2	Distinct elements . . . . .	16
2.2.1	Idealized FM algorithm: freely stored randomness . . . . .	16
2.2.2	A non-idealized algorithm: KMV . . . . .	18
2.2.3	Another algorithm via geometric sampling . . . . .	21
2.3	Quantiles . . . . .	23
2.3.1	q-digest . . . . .	24
2.3.2	MRL . . . . .	27
2.3.3	KLL . . . . .	29
<b>3</b>	<b>Lower Bounds</b>	<b>31</b>
3.1	Compression-based arguments . . . . .	31
3.1.1	Distinct elements . . . . .	31
3.1.2	Quantiles . . . . .	35
3.2	Communication Complexity . . . . .	36
3.2.1	EQUALITY . . . . .	37
3.2.2	DISJOINTNESS . . . . .	38
3.2.3	INDEXING, GAPHAMMING, and DISTINCT ELEMENTS . . . . .	39
<b>4</b>	<b>Linear Sketching</b>	<b>43</b>
4.1	Heavy hitters . . . . .	44
4.1.1	CountMin sketch . . . . .	44
4.1.2	CountSketch . . . . .	47
4.2	Graph sketching . . . . .	49
4.2.1	$k$ -sparse recovery . . . . .	49
4.2.2	SupportFind . . . . .	50
4.2.3	AGM sketch . . . . .	51
4.3	Norm estimation . . . . .	53
4.3.1	AMS sketch . . . . .	53
4.3.2	Indyk's $p$ -stable sketch . . . . .	54

4.3.3	Branching programs and pseudorandom generators . . . . .	56
<b>5</b>	<b>Johnson-Lindenstrauss Transforms</b>	<b>59</b>
5.1	Proof of the Distributional Johnson-Lindenstrauss lemma . . . . .	60
5.2	Lower bound . . . . .	61
5.2.1	Distributional JL . . . . .	62
5.2.2	Optimal JL lower bound . . . . .	62
5.3	Speeding up Johnson-Lindenstrauss transforms . . . . .	65
5.3.1	Sparse Johnson-Lindenstrauss Transform . . . . .	66
5.3.2	Fast Johnson-Lindenstrauss Transform . . . . .	68
5.3.3	Krahmer-Ward theorem . . . . .	70
<b>6</b>	<b>Linear algebra applications</b>	<b>75</b>
6.1	Approximate matrix multiplication . . . . .	75
6.1.1	Sampling approach . . . . .	75
6.1.2	Oblivious linear sketching approach . . . . .	77
6.2	Subspace embeddings . . . . .	79
6.2.1	Given an orthonormal basis . . . . .	80
6.2.2	Leverage score sampling . . . . .	80
6.2.3	Oblivious subspace embeddings . . . . .	82
6.3	Least squares regression . . . . .	84
6.3.1	Sketch-and-solve via subspace embeddings . . . . .	84
6.3.2	Sketch-and-solve via AMM and subspace embeddings . . . . .	85
6.3.3	Accelerating iterative solvers via sketching . . . . .	86
6.4	Approximate low-rank approximation . . . . .	87
6.5	Projection-cost preserving sketches . . . . .	89
6.5.1	$k$ -means clustering . . . . .	93
<b>7</b>	<b>Compressed Sensing</b>	<b>95</b>
7.1	Basis Pursuit . . . . .	97
7.1.1	Obtaining RIP matrices . . . . .	100
7.2	Iterative Hard Thresholding . . . . .	101
<b>8</b>	<b>Suprema of stochastic processes and applications</b>	<b>105</b>
8.1	Methods of bounding gaussian mean width . . . . .	105
8.2	Instance-wise bounds for Johnson-Lindenstrauss . . . . .	108
8.3	Heavy hitters: the BPTree . . . . .	112

# Chapter 1

## Introduction

**Sketching and streaming.** A *sketch*  $C(X)$  of some data set  $X$  with respect to some function  $f$  is a *compression* of  $X$  that allows us to compute, or approximately compute,  $f(X)$  given access only to  $C(X)$ . Sometimes  $f$  has 2 (or multiple) arguments, and for data  $X$  and  $Y$ , we want to compute  $f(X, Y)$  given  $C(X), C(Y)$ . For example, if two servers on a network want to compute some similarity or distance measure on their data, one can simply send the sketch to another (or each to a third party), which reduces network bandwidth compared to sending the entirety of  $X, Y$ .

As a trivial example, consider the case that Alice has a data set  $X$  which is a set of integers, and Bob has a similar data set  $Y$ . They want to compute  $f(X, Y) = \sum_{z \in X \cup Y} z$ . Then each party can let the sketch of their data simply be the sum of all elements in their data set.

When designing *streaming* algorithms, we want to maintain a sketch  $C(X)$  on the fly as  $X$  is updated. In the previous example, if say Alice's data set is being inserted into on the fly then she can of course maintain a sketch by keeping a running sum. The streaming setting appears in many scenarios, such as for example an Internet router monitoring network traffic, or a search engine monitoring a query stream.

### 1.1 Probability Review

We will mainly be dealing with discrete random variables; we consider random variables taking values in some countable subset  $S \subset \mathbb{R}$ . Recall the expectation of  $X$  is defined to be

$$\mathbb{E} X = \sum_{j \in S} j \cdot \mathbb{P}(X = j).$$

We now state a few basic lemmas and facts without proof.

**Lemma 1.1.1** (Linearity of expectation).

$$\mathbb{E}(X + Y) = \mathbb{E} X + \mathbb{E} Y \tag{1.1}$$

**Lemma 1.1.2** (Markov). *If  $X$  is a nonnegative random variable, then*

$$\forall \lambda > 0, \mathbb{P}(X > \lambda) < \frac{\mathbb{E} X}{\lambda}$$

**Lemma 1.1.3** (Chebyshev).

$$\forall \lambda > 0, \mathbb{P}(|X - \mathbb{E} X| > \lambda) < \frac{\mathbb{E}(X - \mathbb{E} X)^2}{\lambda^2} \tag{1.2}$$

*Proof.*  $\mathbb{P}(|X - \mathbb{E}X| > \lambda) = \mathbb{P}((X - \mathbb{E}X)^2 > \lambda^2)$ , and thus the claim follows by Markov's inequality.  $\square$

Rather than the second moment, one can also consider larger moments to obtain:

$$\forall p \geq 1, \forall \lambda > 0, \mathbb{P}(|X - \mathbb{E}X| > \lambda) < \frac{\mathbb{E}|X - \mathbb{E}X|^p}{\lambda^p}. \quad (1.3)$$

By a calculation and picking  $p$  optimally (or by Markov's inequality on the moment-generating function  $e^{tX}$  and appropriately picking  $t$ ), one can also obtain the following ‘‘Chernoff bound’’.

**Theorem 1.1.4** (Chernoff bound). *Suppose  $X_1, \dots, X_n$  are independent random variables with  $X_i \in [0, 1]$ . Let  $X = \sum_i X_i$  and write  $\mu := \mathbb{E}X$ . Then*

$$\forall \lambda > 0, \mathbb{P}(X > (1 + \lambda)\mu) < \left( \frac{e^\lambda}{(1 + \lambda)^{1 + \lambda}} \right)^\mu \quad (\text{upper tail}) \quad (1.4)$$

and also

$$\forall \lambda > 0, \mathbb{P}(X < (1 - \lambda)\mu) < \left( \frac{e^{-\lambda}}{(1 - \lambda)^{1 - \lambda}} \right)^\mu \quad (\text{lower tail}) \quad (1.5)$$

*Proof.* We give the standard proof via bounding the moment-generating function (MGF), and only for the upper tail and only when each  $X_i$  is a Bernoulli random variable with parameter  $p_i$  (equal to 1 with probability  $p_i$  and 0 otherwise); standard references provide proofs of the most general form. The proof for the lower tail is similar. First,

$$\mathbb{P}(X > (1 + \lambda)\mu) = \mathbb{P}(e^{tX} > e^{t(1 + \lambda)\mu})$$

for any  $t \in \mathbb{R}$  since the map  $x \mapsto e^{tx}$  is strictly increasing. As  $e^{tX}$  is a nonnegative random variable, we can apply Markov's inequality.

$$\begin{aligned} \mathbb{P}(e^{tX} > e^{t(1 + \lambda)\mu}) &< e^{-t(1 + \lambda)\mu} \cdot \mathbb{E} e^{tX} && \text{(Markov)} \\ &= e^{-t(1 + \lambda)\mu} \cdot \mathbb{E} e^{t \sum_i X_i} \\ &= e^{-t(1 + \lambda)\mu} \cdot \mathbb{E} \left[ \prod_i e^{tX_i} \right] \\ &= e^{-t(1 + \lambda)\mu} \cdot \prod_i \mathbb{E} e^{tX_i} && \text{(independence of the } X_i) \\ &= e^{-t(1 + \lambda)\mu} \cdot \prod_i (p_i e^t + (1 - p_i)) \\ &= e^{-t(1 + \lambda)\mu} \cdot \prod_i (1 + p_i(e^t - 1)) \\ &\leq e^{-t(1 + \lambda)\mu} \cdot \prod_i e^{p_i(e^t - 1)} && (1 + a \leq e^a) \\ &= e^{-t(1 + \lambda)\mu} \cdot e^{(e^t - 1) \sum_i p_i} \\ &= e^{(e^t - 1 - t - \lambda t)\mu} \\ &= \left( \frac{e^\lambda}{(1 + \lambda)^{1 + \lambda}} \right)^\mu && \text{(set } t = \ln(1 + \lambda)) \end{aligned} \quad (1.6)$$

The lower tail is similar though one writes  $\mathbb{P}(X < (1 - \lambda)\mu) = \mathbb{P}(-X > -(1 - \lambda)\mu) = \mathbb{P}(e^{-tX} > e^{-t(1 - \lambda)\mu})$  then does similar calculations and optimizes choice of  $t$ .  $\square$

**Remark 1.1.5.** The upper tail has two regimes of interest for  $\lambda$ :  $\lambda \ll 1$  and  $\lambda \gg 1$ . Note when  $\lambda < 1$ , we have

$$\begin{aligned} \frac{e^\lambda}{(1+\lambda)^{1+\lambda}} &= \frac{e^\lambda}{e^{(1+\lambda)\ln(1+\lambda)}} \\ &= \frac{e^\lambda}{e^{(1+\lambda)(\lambda - \lambda^2/2 + O(\lambda^3))}} && \text{(Taylor's theorem)} \\ &= e^{-\lambda^2\mu/2 + O(\lambda^3\mu)}. \end{aligned}$$

Rather than use an  $O(\lambda^3)$  bound from Taylor's theorem, one can show a more precise inequality  $\ln(1+\lambda) \geq 2\lambda/(2+\lambda)$  for any  $\lambda \geq 0$  and use it to achieve the final upper bound  $e^{-\lambda^2\mu/3}$ .

The second regime of interest is large  $\lambda$ , i.e.  $\lambda \gg 1$  (specifically  $\lambda > 2e - 1$ ). In this case,  $e/(1+\lambda)$  is less than  $1/2$  and is also  $O(1/\lambda)$  so that the overall upper tail bound is

$$\lambda^{-\Omega(\lambda\mu)}. \quad (1.7)$$

For the lower tail, as we are only ever interested in  $\lambda \leq 1$  (as  $X < 0$  trivially has probability 0 of occurring), logarithmic approximations as for the upper tail can be used to obtain the lower tail bound  $e^{-\lambda^2\mu/2}$ . By combining the upper and lower tails via a union bound,

$$\mathbb{P}(|X - \mu| > \lambda\mu) < e^{-\lambda^2\mu/3} + e^{-\lambda^2\mu/2} < 2e^{-\lambda^2\mu/3}. \quad (1.8)$$

We also have the following, similar ‘‘Hoeffding bound’’, which can also be proven via Markov's inequality applied to the MGF and optimizing choice of  $t$ .

**Theorem 1.1.6** (Hoeffding bound). *Suppose  $X_1, \dots, X_n$  are i.i.d. Bernoulli( $p$ ) random variables for some  $p \in (0, 1)$ . Then for any  $\epsilon \in (0, 1)$ ,*

$$\mathbb{P}\left(\sum_{i=1}^n X_i > (p + \epsilon)n\right) < e^{-2\epsilon^2 n},$$

and

$$\mathbb{P}\left(\sum_{i=1}^n X_i < (p - \epsilon)n\right) < e^{-2\epsilon^2 n}.$$

Another inequality we will make use of is Khintchine's inequality. It essentially says that when  $\sigma_i$  are i.i.d. Rademachers and  $x_i$  are scalars,  $\sum_i \sigma_i x_i$  is what is known as ‘‘subgaussian’’ (i.e. decays at least as fast as a gaussian of some variance).

**Theorem 1.1.7** (Khintchine). *Let  $\sigma_1, \dots, \sigma_n$  be Rademacher random variables (i.e. uniform in  $\{-1, 1\}$ ) and independent, and  $x \in \mathbb{R}^n$  is fixed. Then  $\forall \lambda > 0$ ,  $\mathbb{P}(|\langle \sigma, x \rangle| > \lambda) \leq 2e^{-\lambda^2/2\|x\|_2^2}$ .*

*Proof.* Define  $X = \sum_i \sigma_i x_i$ . Then

$$\begin{aligned} \mathbb{E} e^{tX} &= \prod_i \mathbb{E} e^{t\sigma_i x_i} \\ &= \prod_i \frac{1}{2}(e^{-tx_i} + e^{tx_i}) \\ &= \prod_i \left(1 + \frac{(tx_i)^2}{2!} + \frac{(tx_i)^4}{4!} + \frac{(tx_i)^6}{6!} + \dots\right) && \text{(Taylor expansion)} \end{aligned}$$

$$\begin{aligned}
&\leq \prod_i \left( 1 + \frac{(tx_i)^2}{1! \cdot 2} + \frac{(tx_i)^4}{2! \cdot 2^2} + \frac{(tx_i)^6}{3! \cdot 2^3} + \dots \right) \\
&= \prod_i e^{t^2 x_i^2 / 2} \\
&= e^{t^2 \|x\|_2^2 / 2}
\end{aligned} \tag{1.9}$$

Eq. (1.9) is exactly the MGF of a gaussian random variable  $g$  with mean zero and variance  $\|x\|_2^2$ . To see that this implies the desired tail bound, by Markov's inequality we then have  $\mathbb{P}(X > \lambda) = \mathbb{P}(e^{tX} > e^{t\lambda}) < e^{-t\lambda + t^2 \|x\|_2^2 / 2}$  by Markov's inequality and Eq. (1.9). This is equal to  $e^{-\lambda^2 / (2\|x\|_2^2)}$  by setting  $t = \lambda / \|x\|_2^2$ . Since we are looking at the event  $|X| > \lambda$ , we must also bound  $\mathbb{P}(-X > \lambda) = \mathbb{P}(e^{-tX} > e^{t\lambda}) < e^{-t\lambda + (-t)^2 \|x\|_2^2 / 2}$ , which is the same as above. Thus the claim holds, since  $\mathbb{P}(|X| > \lambda) = \mathbb{P}(X > \lambda) + \mathbb{P}(-X > \lambda)$  as these are disjoint events (or alternatively it is good enough to say the right hand side is an upper bound, which holds via the union bound).  $\square$

We also make use of the following inequality.

**Theorem 1.1.8** (Jensen's inequality). *Let  $X$  be a random variable supported in  $\mathbb{R}$ , and suppose  $\varphi : \mathbb{R} \rightarrow \mathbb{R}$  is convex. Then  $\varphi(\mathbb{E}X) \leq \mathbb{E}\varphi(X)$ .*

The following is a corollary of Jensen's inequality that we frequently use, where  $\|X\|_p$  denotes  $(\mathbb{E}|X|^p)^{1/p}$  is a norm for  $p \geq 1$ .

**Lemma 1.1.9.** *For  $1 \leq p \leq q$ ,  $\|X\|_p \leq \|X\|_q$ .*

*Proof.* Define  $\varphi(x) = |x|^{q/p}$ , which is convex. Then by Theorem 1.1.8 applied to the random variable  $|X|^p$ ,

$$(\mathbb{E}|X|^p)^{q/p} = \varphi(\mathbb{E}|X|^p) \leq \mathbb{E}\varphi(|X|^p) = \mathbb{E}|X|^q.$$

Raising both sides to the  $1/q$ th power yields the result.  $\square$

**Lemma 1.1.10** (Symmetrization / Desymmetrization). *Let  $Z_1, \dots, Z_n$  be independent random variables. Let  $r_1, \dots, r_n$  be independent Rademachers. Then*

$$\left\| \sum_i Z_i - \mathbb{E} \sum_i Z_i \right\|_p \leq 2 \cdot \left\| \sum_i r_i Z_i \right\|_p \text{ (symmetrization inequality)}$$

and

$$(1/2) \cdot \left\| \sum_i r_i (Z_i - \mathbb{E} Z_i) \right\|_p \leq \left\| \sum_i Z_i \right\|_p \text{ (desymmetrization inequality)}.$$

*Proof.* For the first inequality, let  $Y_1, \dots, Y_n$  be independent of the  $Z_i$  but identically distributed to them. Then

$$\begin{aligned}
\left\| \sum_i Z_i - \mathbb{E} \sum_i Z_i \right\|_p &= \left\| \sum_i Z_i - \mathbb{E} \sum_i Y_i \right\|_p \\
&\leq \left\| \sum_i (Z_i - Y_i) \right\|_p && \text{(Jensen)} \\
&= \left\| \sum_i r_i (Z_i - Y_i) \right\|_p && (1.10) \\
&\leq 2 \cdot \left\| \sum_i r_i Z_i \right\|_p && \text{(triangle inequality)}
\end{aligned}$$



Eq. (1.10) follows since the  $X_i - Y_i$  are independent across  $i$  and symmetric.

For the second inequality, let  $Y_i$  be as before. Then

$$\begin{aligned}
\left\| \sum_i r_i (Z_i - \mathbb{E} Z_i) \right\|_p &= \left\| \mathbb{E}_Y \sum_i r_i (Z_i - Y_i) \right\|_p \\
&\leq \left\| \sum_i r_i (Z_i - Y_i) \right\|_p && \text{(Jensen)} \\
&= \left\| \sum_i (Z_i - Y_i) \right\|_p \\
&\leq 2 \cdot \left\| \sum_i Z_i \right\|_p && \text{(triangle inequality)}
\end{aligned}$$

□

**Lemma 1.1.11** (Decoupling [dlPnG99]). *Let  $x_1, \dots, x_n$  be independent and mean zero, and  $x'_1, \dots, x'_n$  identically distributed as the  $x_i$  and independent of them. Then for any  $(a_{i,j})$  and for all  $p \geq 1$*

$$\left\| \sum_{i \neq j} a_{i,j} x_i x_j \right\|_p \leq 4 \left\| \sum_{i,j} a_{i,j} x_i x'_j \right\|_p$$

*Proof.* Let  $\eta_1, \dots, \eta_n$  be independent Bernoulli random variables each of expectation 1/2. Then

$$\begin{aligned}
\left\| \sum_{i \neq j} a_{i,j} x_i x_j \right\|_p &= 4 \cdot \left\| \mathbb{E}_\eta \sum_{i \neq j} a_{i,j} x_i x_j |\eta_i| |1 - \eta_j| \right\|_p \\
&\leq 4 \cdot \left\| \sum_{i \neq j} a_{i,j} x_i x_j \eta_i (1 - \eta_j) \right\|_p && \text{(Jensen)} \quad (1.11)
\end{aligned}$$

Hence there must be some fixed vector  $\eta' \in \{0, 1\}^n$  which achieves

$$\left\| \sum_{i \neq j} a_{i,j} x_i x_j \eta_i (1 - \eta_j) \right\|_p \leq \left\| \sum_{i \in S} \sum_{j \notin S} a_{i,j} x_i x_j \right\|_p$$

where  $S = \{i : \eta'_i = 1\}$ . Let  $x_S$  denote the  $|S|$ -dimensional vector corresponding to the  $x_i$  for  $i \in S$ . Then

$$\begin{aligned}
\left\| \sum_{i \in S} \sum_{j \notin S} a_{i,j} x_i x_j \right\|_p &= \left\| \sum_{i \in S} \sum_{j \notin S} a_{i,j} x_i x'_j \right\|_p \\
&= \left\| \mathbb{E}_{x_S} \mathbb{E}_{x'_S} \sum_{i,j} a_{i,j} x_i x'_j \right\|_p \quad (\mathbb{E} x_i = \mathbb{E} x'_j = 0) \\
&\leq \left\| \sum_{i,j} a_{i,j} x_i x'_j \right\|_p && \text{(Jensen)}
\end{aligned}$$

□

The following proof of the Hanson-Wright was shared to me by Sjoerd Dirksen (personal communication). A newer proof which we do not cover here, using more modern tools that allow extension to subgaussian variables and not just Rademachers, is given in [RV13]. Recall the Frobenius norm is defined by  $\|A\|_F := (\sum_{i,j} A_{i,j}^2)^{1/2}$ , and the operator norm by  $\|A\| = \sup_{\|x\|_2 = \|y\|_2 = 1} x^\top A y$ .

**Theorem 1.1.12** (Hanson-Wright inequality [HW71]). For  $\sigma_1, \dots, \sigma_n$  independent Rademachers and  $A \in \mathbb{R}^{n \times n}$ , for all  $p \geq 1$

$$\|\sigma^\top A \sigma - \mathbb{E} \sigma^\top A \sigma\|_p \lesssim \sqrt{p} \cdot \|A\|_F + p \cdot \|A\|.$$

*Proof.* Without loss of generality we assume in this proof that  $p \geq 2$  (so that  $p/2 \geq 1$ ). Then

$$\|\sigma^\top A \sigma - \mathbb{E} \sigma^\top A \sigma\|_p \lesssim \|\sigma^\top A \sigma'\|_p \quad (\text{Lemma 1.1.11}) \quad (1.12)$$

$$\lesssim \sqrt{p} \cdot \|\|Ax\|_2\|_p \quad (\text{Khintchine}) \quad (1.13)$$

$$= \sqrt{p} \cdot \|\|Ax\|_2^2\|_{p/2}^{1/2} \quad (1.14)$$

$$\leq \sqrt{p} \cdot \|\|Ax\|_2^2\|_p^{1/2}$$

$$\leq \sqrt{p} \cdot (\|A\|_F^2 + \|\|Ax\|_2^2 - \mathbb{E} \|Ax\|_2^2\|_p)^{1/2} \quad (\text{triangle inequality})$$

$$\leq \sqrt{p} \cdot \|A\|_F + \sqrt{p} \cdot \|\|Ax\|_2^2 - \mathbb{E} \|Ax\|_2^2\|_p^{1/2}$$

$$\lesssim \sqrt{p} \cdot \|A\|_F + \sqrt{p} \cdot \|x^\top A^\top A x'\|_p^{1/2} \quad (\text{Lemma 1.1.11})$$

$$\lesssim \sqrt{p} \cdot \|A\|_F + p^{3/4} \cdot \|\|A^\top A x\|_2\|_p^{1/2} \quad (\text{Khintchine})$$

$$\lesssim \sqrt{p} \cdot \|A\|_F + p^{3/4} \cdot \|A\|^{1/2} \cdot \|\|Ax\|_2\|_p^{1/2} \quad (1.15)$$

Writing  $E = \|\|Ax\|_2\|_p^{1/2}$  and comparing Eq. (1.13) and Eq. (1.15), we see that for some constant  $C > 0$ ,

$$E^2 - Cp^{1/4}\|A\|^{1/2}E - C\|A\|_F \leq 0.$$

Thus  $E$  must be smaller than the larger root of the above quadratic equation, implying our desired upper bound on  $E^2$ .  $\square$

**Remark 1.1.13.** The ‘‘square root trick’’ in the proof of the Hanson-Wright inequality above is quite handy and can be used to prove several moment inequalities (for example, it can be used to prove Bernstein inequality). As far as I am aware, the trick was first used in a work of Rudelson [Rud99] on operator norms of certain random matrices.

**Remark 1.1.14.** We could have upper bounded Eq. (1.14) by

$$\sqrt{p} \cdot \|A\|_F + \sqrt{p} \cdot \|\|Ax\|_2^2 - \mathbb{E} \|Ax\|_2^2\|_{p/2}^{1/2}$$

by the triangle inequality. Now notice we have bounded the  $p$ th central moment of a symmetric quadratic form Eq. (1.12) by the  $p/2$ th moment also of a symmetric quadratic form. Writing  $p = 2^k$ , this observation leads to a proof by induction on  $k$ , which was the approach used in [DKN10].

We have stated a moment version of the Hanson-Wright inequality, but we show now this is equivalent to a tail bound. Below we prove a lemma which lets us freely obtain tail bounds from moment bounds and vice versa (often we prove a moment bound and later invoke a tail bound, or vice versa, without even mentioning any justification).

**Lemma 1.1.15.** Let  $Z$  be a scalar random variable. Consider the following statements:

(1a) There exists  $\sigma > 0$  s.t.  $\forall p \geq 1$ ,  $\|Z\|_p \leq C_1 \sigma \sqrt{p}$ .

(1b) There exists  $\sigma > 0$  s.t.  $\forall \lambda > 0$ ,  $\mathbb{P}(|Z| > \lambda) \leq C_2 e^{-C_2' \lambda^2 / \sigma^2}$ .

(2a) There exists  $K > 0$  s.t.  $\forall p \geq 1$ ,  $\|Z\|_p \leq C_3 K p$ .

(2b) There exists  $K > 0$  s.t.  $\forall \lambda > 0$ ,  $\mathbb{P}(|Z| > \lambda) \leq C_4 e^{-C'_4 \lambda/K}$ .

(3a) There exist  $\sigma, K > 0$  s.t.  $\forall p \geq 1$ ,  $\|Z\|_p \leq C_5(\sigma\sqrt{p} + Kp)$ .

(3b) There exist  $\sigma, K > 0$  s.t.  $\forall \lambda > 0$ ,  $\mathbb{P}(|Z| > \lambda) \leq C_6(e^{-C'_6 \lambda^2/\sigma^2} + e^{-C'_6 \lambda/K})$ .

Then 1a is equivalent to 1b, 2a is equivalent to 2b, and 3a is equivalent to 3b, where the constants  $C_i, C'_i$  in each case change by at most some absolute constant factor.

*Proof.* We will show only that 1a is equivalent to 1b; the other cases are argued identically.

To show that 1a implies 1b, by Markov's inequality

$$\mathbb{P}(Z > \lambda) \leq \lambda^{-p} \cdot \mathbb{E}|Z|^p \leq \left( \frac{C_1^2 \sigma^2}{\lambda^2 p} \right)^{p/2}.$$

Statement 1b follows by choosing  $p = \max\{1, 2C_1^2 \lambda^2/\sigma^2\}$ .

To show that 1b implies 1a, by integration by parts we have

$$\mathbb{E}|Z|^p = \int_0^\infty p x^{p-1} \mathbb{P}(|Z| > \lambda) d\lambda \leq 2C_2 p \cdot \int_0^\infty p x^{p-1} \cdot e^{-C'_2 \lambda^2/\sigma^2} d\lambda.$$

The integral on the right hand side is exactly the  $p$ th moment of a gaussian random variable with mean zero and variance  $\sigma'^2 = \sigma^2/(2C'_2)$ . Statement 1a then follows since such a gaussian has  $p$ -norm  $\Theta(\sigma' \sqrt{p})$ .  $\square$

**Corollary 1.1.16.** For  $\sigma_1, \dots, \sigma_n$  independent Rademachers and  $A \in \mathbb{R}^{n \times n}$ , for all  $\lambda > 0$

$$\mathbb{P}(|\sigma^\top A \sigma - \mathbb{E} \sigma^\top A \sigma| > \lambda) \lesssim e^{-C\lambda^2/\|A\|_F^2} + e^{-C\lambda/\|A\|}.$$

*Proof.* Combine [Theorem 1.1.12](#) with item (3b) of [Lemma 1.1.15](#).  $\square$

As mentioned in [Remark 1.1.13](#), we now show that the ‘‘square root trick’’ can also be used to prove Bernstein's inequality.

**Theorem 1.1.17** (Bernstein's inequality; moment form). *Let  $X_1, \dots, X_n$  be independent, each bounded in magnitude by  $K$  almost surely. Write  $X := \sum_{i=1}^n X_i$  and define  $\sigma^2 := \text{Var}[X]$ . Then*

$$\forall p \geq 1, \|X - \mathbb{E} X\|_p \lesssim \sigma\sqrt{p} + Kp.$$

*Proof.* Define  $Z = X_i - \mathbb{E} X_i$ . Then

$$\begin{aligned} \|Z\|_p &= \|Z - \mathbb{E} Z\|_p \\ &\leq 2 \left\| \sum_i \sigma_i Z_i \right\|_p && \text{(Lemma 1.1.10)} \\ &\leq 2\sqrt{p} \left\| \left( \sum_i Z_i^2 \right)^{1/2} \right\|_p && \text{(Khintchine)} \\ &= 2\sqrt{p} \left\| \sum_i Z_i^2 \right\|_{p/2}^{1/2} && (1.16) \\ &\leq 2\sqrt{p} \left\| \sum_i Z_i^2 \right\|_p^{1/2} \\ &\leq 2\sqrt{p} \left( \mathbb{E} \sum_i Z_i^2 \right)^{1/2} + 2\sqrt{p} \left\| \sum_i Z_i^2 - \mathbb{E} \sum_i Z_i^2 \right\|_p^{1/2} && \text{(triangle inequality)} \end{aligned}$$

$$\begin{aligned}
&= 2\sigma\sqrt{p} + 2\sqrt{p}\left\|\sum_i Z_i^2 - \mathbb{E}\sum_i Z_i^2\right\|_p^{1/2} \\
&\leq 2\sigma\sqrt{p} + 4\sqrt{p}\left\|\sum_i \sigma_i Z_i^2\right\|_p^{1/2} && \text{(Lemma 1.1.10)} \\
&\leq 2\sigma\sqrt{p} + 4p^{3/4}\left\|\left(\sum_i Z_i^4\right)^{1/2}\right\|_p^{1/2} && \text{(Khintchine)} \\
&\leq 2\sigma\sqrt{p} + 4p^{3/4}\sqrt{K}\left\|\left(\sum_i Z_i^2\right)^{1/2}\right\|_p^{1/2} \\
&= 2\sigma\sqrt{p} + 4p^{3/4}\sqrt{K}\left\|\sum_i Z_i^2\right\|_{p/2}^{1/4} \tag{1.17}
\end{aligned}$$

Define  $E := \left\|\sum_i Z_i^2\right\|_{p/2}^{1/4}$ . Then comparing Eqs. (1.16) and (1.17), we find

$$2\sqrt{p}E^2 \leq 2\sigma\sqrt{p} + 4p^{3/4}\sqrt{K}E.$$

Rearranging gives

$$E^2 - 2p^{1/4}\sqrt{K}E - \sigma \leq 0.$$

Thus  $E$  must be smaller than the larger root of the associated quadratic, i.e.  $E \leq (2p^{1/4}\sqrt{K} + \sqrt{4K\sqrt{p} + 4\sigma})/2 \lesssim p^{1/4}\sqrt{K} + \sqrt{\sigma}$ . The claim then follows since we have  $\|X - \mathbb{E}X\|_p \leq 2\sqrt{p}E^2$ .  $\square$

**Corollary 1.1.18** (Bernstein's inequality; tail form). *Let  $X_1, \dots, X_n$  be independent, each bounded in magnitude by  $K$  almost surely. Write  $X := \sum_{i=1}^n X_i$  and define  $\sigma^2 := \text{Var}[X]$ . Then*

$$\forall \lambda > 0, \mathbb{P}(|X - \mathbb{E}X| > \lambda) \lesssim e^{-C\lambda^2/\sigma^2} + e^{-C\lambda/K}$$

for some universal constant  $C > 0$ .

*Proof.* Apply the equivalence of (3a) and (3b) of Lemma 1.1.15 to Theorem 1.1.17.  $\square$

## Chapter 2

# Counting Problems

### 2.1 Approximate counting

In the following, we discuss a problem first studied in [Mor78].

**Problem.** Our algorithm must monitor a sequence of events, then at any given time output (an estimate of) the number of events thus far. More formally, this is a data structure maintaining a single integer  $n$  and supporting the following two operations:

- **update():** increment  $n$  by 1
- **query():** output (an estimate of)  $n$

Before any operations are performed, it is assumed that  $n$  starts at 0. Of course a trivial algorithm maintains  $n$  using  $O(\log n)$  bits of memory (a counter). Our goal is to use much less space than this. It is not too hard to prove that it is impossible to solve this problem exactly using  $o(\log n)$  bits of space. Thus we would like to answer **query()** with some estimate  $\tilde{n}$  of  $n$  satisfying

$$\mathbb{P}(|\tilde{n} - n| > \varepsilon n) < \delta, \tag{2.1}$$

for some  $0 < \varepsilon, \delta < 1$  that are given to the algorithm up front.

The algorithm of Morris provides such an estimator for some  $\varepsilon, \delta$  that we will analyze shortly. The algorithm works as follows:

1. Initialize  $X \leftarrow 0$ .
2. For each update, increment  $X$  with probability  $\frac{1}{2^X}$ .
3. For a query, output  $\tilde{n} = 2^X - 1$ .

Intuitively, the variable  $X$  is attempting to store a value that is  $\approx \log_2 n$ . Before giving a rigorous analysis in [Subsection 2.1.1](#), we first give a probability review.

#### 2.1.1 Analysis of Morris' algorithm

Let  $X_n$  denote  $X$  in Morris' algorithm after  $n$  updates.

**Claim 2.1.1.**

$$\mathbb{E} 2^{X_n} = n + 1.$$

*Proof.* We prove by induction on  $n$ . The base case is clear, so we now show the inductive step.

We have

$$\begin{aligned}
\mathbb{E} 2^{X_{n+1}} &= \sum_{j=0}^{\infty} \mathbb{P}(X_n = j) \cdot \mathbb{E}(2^{X_{n+1}} | X_n = j) \\
&= \sum_{j=0}^{\infty} \mathbb{P}(X_n = j) \cdot \left(2^j \left(1 - \frac{1}{2^j}\right) + \frac{1}{2^j} \cdot 2^{j+1}\right) \\
&= \sum_{j=0}^{\infty} \mathbb{P}(X_n = j) 2^j + \sum_j \mathbb{P}(X_n = j) \\
&= \mathbb{E} 2^{X_n} + 1 \\
&= (n + 1) + 1
\end{aligned} \tag{2.2}$$

□

It is now clear why we output our estimate of  $n$  as  $\tilde{n} = 2^X - 1$ : it is an unbiased estimator of  $n$ . In order to show [Eq. \(2.1\)](#) however, we will also control on the variance of our estimator. This is because, by Chebyshev's inequality,

$$\mathbb{P}(|\tilde{n} - n| > \varepsilon n) < \frac{1}{\varepsilon^2 n^2} \cdot \mathbb{E}(\tilde{n} - n)^2 = \frac{1}{\varepsilon^2 n^2} \mathbb{E}(2^{2X} - 1 - n)^2. \tag{2.3}$$

When we expand the above square, we find that we need to control  $\mathbb{E} 2^{2X_n}$ . The proof of the following claim is by induction, similar to that of [Claim 2.1.1](#).

**Claim 2.1.2.**

$$\mathbb{E}(2^{2X_n}) = \frac{3}{2}n^2 + \frac{3}{2}n + 1. \tag{2.4}$$

This implies  $\mathbb{E}(\tilde{n} - n)^2 = (1/2)n^2 - (1/2)n - 1 < (1/2)n^2$ , and thus

$$\mathbb{P}(|\tilde{n} - n| > \varepsilon n) < \frac{1}{\varepsilon^2 n^2} \cdot \frac{n^2}{2} = \frac{1}{2\varepsilon^2}, \tag{2.5}$$

which is not particularly meaningful since the right hand side is only better than 1/2 failure probability when  $\varepsilon \geq 1$  (which means the estimator may very well always be 0!).

### 2.1.2 Morris+

To decrease the failure probability of Morris' basic algorithm, we instantiate  $s$  independent copies of Morris' algorithm and average their outputs. That is, we obtain independent estimators  $\tilde{n}_1, \dots, \tilde{n}_s$  from independent instantiations of Morris' algorithm, and our output to a query is

$$\tilde{n} = \frac{1}{s} \sum_{i=1}^s \tilde{n}_i$$

Since each  $\tilde{n}_i$  is an unbiased estimator of  $n$ , so is their average. Furthermore, since variances of independent random variables add, and multiplying a random variable by some constant  $c = 1/s$  causes the variance to be multiplied by  $c^2$ , the right hand side of [\(Eq. \(2.5\)\)](#) becomes

$$\mathbb{P}(|\tilde{n} - n| > \varepsilon n) < \frac{1}{2s\varepsilon^2} < \delta$$

for  $s > 1/(2\varepsilon^2\delta) = \Theta(1/(\varepsilon^2\delta))$ .

### 2.1.3 Morris++

It turns out there is a simple technique (which we will see often) to reduce the dependence on the failure probability  $\delta$  from  $1/\delta$  to  $\log(1/\delta)$ . The technique is as follows.

We run  $t$  instantiations of Morris+, each with failure probability  $\frac{1}{3}$ . Thus, for each one,  $s = \Theta(1/\varepsilon^2)$ . We then output the median estimate from all the  $s$  Morris+ instantiations. Note that the expected number of Morris+ instantiations that succeed is at least  $2t/3$ . For the median to be a bad estimate, at most half the Morris+ instantiations can succeed, implying the number of succeeding instantiations deviated from its expectation by at least  $t/6$ . Define

$$Y_i = \begin{cases} 1, & \text{if the } i\text{th Morris+ instantiation succeeds.} \\ 0, & \text{otherwise.} \end{cases} \quad (2.6)$$

Then by the Hoeffding bound,

$$\mathbb{P}\left(\sum_i Y_i \leq \frac{t}{2}\right) \leq \mathbb{P}\left(\sum_i Y_i - \mathbb{E}\sum_i Y_i < -t/6\right) < e^{-2(1/6)^2 t} \leq \delta \quad (2.7)$$

for  $t \geq \lceil 18 \ln(1/\delta) \rceil$ .

**Overall space complexity.** When we unravel Morris++, it is running a total of  $st = \Theta(\lg(1/\delta)/\varepsilon^2)$  instantiations of the basic Morris algorithm. Now note that once any given Morris counter  $X$  reaches the value  $\lg(stn/\delta)$ , the probability that it is incremented at any given moment is at most  $\delta/(nst)$ . Thus the probability that it is incremented at all in the next  $n$  increments is at most  $\delta/(st)$ . Thus by a union bound, with probability  $1 - \delta$  none of the  $st$  basic Morris instantiations ever stores a value larger than  $\lg(stn/\delta)$ , which takes  $O(\lg \lg(stn/\delta))$  bits. Thus the total space complexity is, with probability  $1 - \delta$ , at most

$$O(\varepsilon^{-2} \lg(1/\delta) (\lg \lg(n/(\varepsilon\delta)))).$$

In particular, for constant  $\varepsilon, \delta$  (say each  $1/100$ ), the total space complexity is  $O(\lg \lg n)$  with constant probability. This is exponentially better than the  $\log n$  space achieved by storing a counter!

**Remark 2.1.3.** As we will continue to see, designing some random process and associated unbiased estimator with bounded variance for some desired statistic that can be maintained and computed in small space is a common strategy. One can then take the median of means of several copies of this basic structure to then obtain  $1 + \varepsilon$ -approximation with success probability  $1 - \delta$ . In the particular case of approximate counting, this led to a space blow-up of  $\Theta(\log(1/\delta)/\varepsilon^2)$  though it is possible to do better by a more tailored approach. Specifically, Morris instead suggests adjusting the increment probability from  $1/2^X$  to  $1/(1 + \alpha)^X$  for small  $\alpha$  then estimating  $n$  as  $\tilde{n} := ((1 + \alpha)^X - 1)/\alpha$ . This makes intuitive sense: if we increment with probability  $1.0^X$ , then we have a simple deterministic counter, which has zero variance but poor memory performance. Meanwhile if we increment with probability  $0.5^X$  then the memory usage is reduced at the cost of higher variance. One may then intuit that if incrementing with probability  $1/(1 + \alpha)^X$ , space usage increases while variance decreases as  $\alpha \downarrow 0$ , which is indeed the case. By bounding variance and applying Chebyshev's inequality, it is possible to show that  $\alpha = \Theta(\varepsilon^2 \delta)$  leads to  $(1 + \varepsilon)$ -approximation with probability  $1 - \delta$  while using space  $O(\log(1/\varepsilon) + \log \log n + \log(1/\delta))$  with high probability [Mor78, Fla85]. In fact, it is even possible to improve the analysis to show  $O(\log(1/\varepsilon) + \log \log n + \log \log(1/\delta))$  space suffices via a different argument, and to show that this bound is optimal [NY20].

## 2.2 Distinct elements

We next consider another counting problem: the *count distinct* or *distinct elements problem*, also known as the  $F_0$  problem, defined as follows. We are given a stream of integers  $i_1, \dots, i_m \in [n]$  where  $[n]$  denotes the set  $\{1, 2, \dots, n\}$ . We would like to output the number of *distinct* elements seen in the stream. As with Morris' approximate counting algorithm, our goal will be to minimize our space consumption.<sup>1</sup>

There are two straightforward solutions as follows:

1. **Solution 1:** keep a bit array of length  $n$ , initialized to all zeroes. Set the  $i$ th bit to 1 whenever  $i$  is seen in the stream ( $n$  bits of memory).
2. **Solution 2:** Store the whole stream in memory explicitly ( $m \lceil \log_2 n \rceil$  bits of memory).

We can thus solve the problem exactly using  $\min\{n, m \lceil \log_2 n \rceil\}$  bits of memory. In order to reduce the space, we will instead settle for computing some value  $\tilde{t}$  s.t.  $\mathbb{P}(|t - \tilde{t}| > \varepsilon t) < \delta$ , where  $t$  denotes the number of distinct elements in the stream. The first work to show that this is possible using small memory (assuming oracle access to certain random hash functions) is that of Flajolet and Martin (FM) [FM85].

### 2.2.1 Idealized FM algorithm: freely stored randomness

We first discuss the following idealized algorithm:

1. Pick a random function  $h : [n] \rightarrow [0, 1]$
2. Maintain counter  $X = \min_{i \in \text{stream}} h(i)$
3. Output  $1/X - 1$

Note this algorithm really is idealized, since we cannot afford to store a truly random such function  $h$  in  $o(n)$  bits (first, because there are  $n$  independent random variables  $(h(i))_{i=1}^n$ , and second because its outputs are real numbers).

**Intuition.** The value  $X$  stored by the algorithm is a random variable that is the minimum of  $t$  i.i.d  $Unif(0, 1)$  random variables. One can then compute the expectation of what the minimum is, as a function of  $t$ , and invert that function in hopes of obtaining (a good approximation to)  $t$ .

**Claim 2.2.1.**  $\mathbb{E} X = \frac{1}{t+1}$ .

*Proof.*

$$\begin{aligned} \mathbb{E} X &= \int_0^\infty \mathbb{P}(X > \lambda) d\lambda \\ &= \int_0^\infty \mathbb{P}(\forall i \in \text{stream}, h(i) > \lambda) d\lambda \\ &= \int_0^\infty \prod_{r=1}^t \mathbb{P}(h(i_r) > \lambda) d\lambda \end{aligned}$$

<sup>1</sup>The reason for the name " $F_0$ " is that if one defines a histogram  $x \in \mathbb{R}^n$  with  $x_i$  being the number of occurrences of  $i$  in the stream, then one can define the  $p$ th moment  $F_p$  as  $\sum_i x_i^p$ . The number of distinct elements is then the limit of the  $p$ th moment for  $p \rightarrow 0$ .



$$\begin{aligned}
&= \int_0^1 (1 - \lambda)^t d\lambda \\
&= \frac{1}{t+1}
\end{aligned}$$

□

We will also need the following claim in order to execute Chebyshev's inequality to bound the failure probability in our final algorithm.

**Claim 2.2.2.**  $\mathbb{E} X^2 = \frac{2}{(t+1)(t+2)}$

*Proof.*

$$\begin{aligned}
\mathbb{E} X^2 &= \int_0^1 \mathbb{P}(X^2 > \lambda) d\lambda \\
&= \int_0^1 \mathbb{P}(X > \sqrt{\lambda}) d\lambda \\
&= \int_0^1 (1 - \sqrt{\lambda})^t d\lambda \\
&= 2 \int_0^1 u^t (1 - u) du && \text{(substitution } u = 1 - \sqrt{\lambda}\text{)} \\
&= \frac{2}{(t+1)(t+2)}
\end{aligned}$$

□

This gives  $\text{Var}[X] = \mathbb{E} X^2 - (\mathbb{E} X)^2 = \frac{t}{(t+1)^2(t+2)}$ , or the simpler  $\text{Var}[X] < (\mathbb{E} X)^2 = \frac{1}{(t+1)^2}$ .

**FM+.** To obtain an algorithm providing a randomized approximate guarantee, just as with Morris+ we form an algorithm FM+ which averages together the outputs from  $s$  independent instantiations of the basic FM algorithm.

1. Instantiate  $s = \lceil 1/(\epsilon^2 \delta) \rceil$  FMs independently,  $\text{FM}_1, \dots, \text{FM}_s$ .
2. Let  $X_i$  be the output of  $\text{FM}_i$ .
3. Upon a query, output  $1/Z - 1$ , where  $Z = \frac{1}{s} \sum_i X_i$ .

We have that  $\mathbb{E}(Z) = \frac{1}{t+1}$ , and  $\text{Var}(Z) = \frac{1}{s} \frac{t}{(t+1)^2(t+2)} < \frac{1}{s(t+1)^2}$ .

**Claim 2.2.3.**  $\mathbb{P}(|Z - \frac{1}{t+1}| > \frac{\epsilon}{t+1}) < \delta$

*Proof.* We apply Chebyshev's inequality.

$$\mathbb{P}(|Z - \frac{1}{t+1}| > \frac{\epsilon}{t+1}) < \frac{(t+1)^2}{\epsilon^2} \frac{1}{s(t+1)^2} = \eta$$

□

**Claim 2.2.4.**  $\mathbb{P}(|(\frac{1}{Z} - 1) - t| > O(\epsilon)t) < \eta$

*Proof.* If  $t = 0$  (the empty stream), the claim is trivial. Otherwise, by the previous claim, with probability  $1 - \eta$  we have

$$\frac{1}{(1 \pm \epsilon)^{\frac{1}{t+1}}} - 1 = (1 \pm O(\epsilon))(t+1) - 1 = (1 \pm O(\epsilon))t \pm O(\epsilon) = (1 \pm O(\epsilon))t$$

□

**FM++.** To obtain our final algorithm, again as with Morris++ we take the median output from multiple independent instantiations of FM+.

1. Instantiate  $q = \lceil 18 \ln(1/\delta) \rceil$  independent copies of FM+ with  $\eta = 1/3$ .
2. Output the median  $\hat{t}$  of  $\{1/Z_j - 1\}_{j=1}^q$  where  $Z_j$  is the output of the  $j$ th copy of FM+.

**Claim 2.2.5.**  $\mathbb{P}(|\hat{t} - t| > \epsilon t) < \delta$

*Proof.* Let

$$Y_j = \begin{cases} 1 & \text{if } |(1/Z_j - 1) - t| \leq \epsilon t \\ 0 & \text{else} \end{cases}$$

and put  $Y = \sum_{j=1}^q Y_j$ . We have  $\mathbb{E}Y > 2q/3$  by our choice of  $\eta$ . The probability we seek to bound is equivalent to the probability that the median fails, i.e. at least half of the FM+ estimates have  $Y_j = 0$ . In other words,

$$\sum_{j=1}^q Y_j \leq q/2$$

We then get that

$$\mathbb{P}\left(\sum Y_j < q/2\right) \leq \mathbb{P}(Y - \mathbb{E}Y < -q/6) \quad (2.8)$$

Then by a Hoeffding bound, the above is at most

$$e^{-2(\frac{1}{6})^2 q} \leq \delta$$

as desired. □

The final space required, ignoring  $h$ , is that required to store  $O(\frac{\lg(1/\delta)}{\epsilon^2})$  real numbers since  $O(\lg(1/\delta))$  copies of FM+ are instantiated, each averaging  $O(1/\epsilon^2)$  copies of FM. Each single FM just stores a single number (the minimum hash ever seen).

### 2.2.2 A non-idealized algorithm: KMV

We next describe a modified algorithm for  $F_0$ -estimation which does not assume access to a truly random hash function. Before we continue, we first discuss  $k$ -wise independent hash families.

**$k$ -wise independent hash families.**

**Definition 2.2.6.** A family  $\mathcal{H}$  of functions mapping  $[a]$  to  $[b]$  is  $k$ -wise independent if  $\forall j_1, \dots, j_k \in [b]$  and  $\forall$  distinct  $i_1, \dots, i_k \in [a]$ ,

$$\mathbb{P}_{h \in \mathcal{H}}(h(i_1) = j_1 \wedge \dots \wedge h(i_k) = j_k) = 1/b^k$$

**Example.** The set  $\mathcal{H}$  of all functions  $[a] \rightarrow [b]$  is  $k$ -wise independent for every  $k$ .  $|\mathcal{H}| = b^a$  so  $h \in \mathcal{H}$  is representable in  $a \lg b$  bits.

**Example.** Let  $a = b = q$  for  $q = p^r$  a prime power and define  $\mathcal{H}_{poly}$  to be the set of all degree  $\leq k - 1$  polynomials with coefficients in  $\mathbb{F}_q$ , the finite field of order  $q$ .  $|\mathcal{H}_{poly}| = q^k$  so  $h \in \mathcal{H}$  is representable in  $k \lg p = k \lg a$  bits.

**Claim 2.2.7** ([WC79]).  $\mathcal{H}_{poly}$  is  $k$ -wise independent.

*Proof.* Given  $((i_r, j_r))_{r=1}^k$ , there is exactly one degree at most  $k - 1$  polynomial  $h$  over  $\mathbb{F}_q$  with  $h(i_r) = j_r$  for  $r = 1, \dots, k$ , via interpolation. Thus the probability

$$\mathbb{P}_{h \in \mathcal{H}}(h(i_1) = j_1 \wedge \dots \wedge h(i_k) = j_k)$$

exactly equals  $1/|\mathcal{H}_{poly}| = 1/p^k = 1/b^k$ . □

We now present an algorithm of [BYJK<sup>+</sup>02], later known as the “KMV algorithm” (for “ $k$  minimum values”). We will assume  $1/\varepsilon^2 < n$ , since we will be shooting for a space bound that is at least  $1/\varepsilon^2$ , and there is always a trivial solution for exact  $F_0$  computation using  $n$  bits. We also assume, without loss of generality, that  $\varepsilon < 1/2$ .

The algorithm is quite similar to the idealized FM algorithm, but rather than maintain only the *smallest* hash evaluation, we maintain the  $k$  smallest hash evaluations for some appropriately chosen  $k \in \Theta(1/\varepsilon^2)$ . The intuition is that the smallest hash evaluation leads us to an estimator with high variance (all it takes is for one item to hash to something really tiny, which will throw off our estimator). Meanwhile, if we consider the  $k$ th smallest hash evaluation for some large  $k$ , then the variance should be smaller since a single item (or small number of items) cannot throw off our statistic by having a wildly small or big hash value. This makes some intuitive sense since the median ( $k = t/2$ ) is the most robust order statistic to outliers. The tradeoff of course is that to keep track of the  $k$ th smallest hash value we will use space that grows with  $k$ , as we will keep track of all  $k$  bottom hash values seen.

Specifically, the KMV algorithm chooses a hash function  $h : [n] \rightarrow [M]$  for  $M = n^3$  from a 2-wise independent hash family (the idea here is to discretize  $[0, 1]$  into multiples of  $1/M$ ). We pick  $k = \lceil 24/\varepsilon^2 \rceil$ . We keep track in memory of the  $k$  smallest hash evaluations. If at the time of the query we have seen less than  $k$  distinct hash values, then we just output the number of distinct hash values seen. Otherwise, if  $X$  is the  $k$ th smallest then we output our estimate of  $t = F_0$  as  $\tilde{t} = kM/X$ .

For some intuition: note if we had  $t \geq k$  independent hash values in  $[0, 1]$ , we expect the  $k$ th smallest value  $v$  to be  $k/(t + 1)$  (namely, we expect  $k$  equally spaced values between 0 and 1). Thus a reasonable estimate for  $t$  would be  $k/v - 1$ . If  $t \leq k$  then our output is exactly correct as long as the distinct elements are perfectly hashed, and indeed this happens with high probability since all of  $[n]$  is perfectly hashed with large probability by our large choice of  $M$ . Otherwise, for  $t \geq k \gg 1/\varepsilon^2$ , we expect  $k/v > 1/\varepsilon^2$ , and thus neglecting to subtract the 1 and simply outputting  $k/v$  gives the same answer up to a  $1 + \varepsilon$  factor. Since in our actual algorithm we discretize  $[0, 1]$  into multiples of  $1/M$ , our value  $X$  is actually representing  $Mv$ , and thus we output  $kM/X$ .

We now provide a more formal analysis. Note that our hash function  $h$  is only 2-wise independent! We would like to say that with good probability,

$$(1 - \varepsilon)t \leq \tilde{t} \leq (1 + \varepsilon)t.$$

We consider the two bad events that  $\tilde{t}$  is too big or too small, and show that each happens with probability at most  $1/6$ , and thus the probability that either happens is at most  $1/3$  by the union bound.

First let us consider the case that  $\tilde{t} > (1 + \varepsilon)t$ . Since  $\tilde{t} = kM/X$ , i.e.  $X = kM/\tilde{t}$ , this can only happen if at least  $k$  distinct indices in the stream hashed to a value smaller than  $kM/((1 + \varepsilon)t)$ . Let  $Y_i$  be an indicator random variable for the event that the  $i$ th distinct integer in the stream hashed to a value below  $kM/((1 + \varepsilon)t)$ , and let  $Y$  denote  $\sum_{i=1}^t Y_i$ . Then  $\mathbb{E} Y_i < k/((1 + \varepsilon)t)$  and thus

$$\mathbb{E} Y < k/(1 + \varepsilon).$$

We also have  $\text{Var}[Y_i] < \mathbb{E} Y_i^2 = \mathbb{E} Y_i < k/((1 + \varepsilon)t)$ , and thus

$$\text{Var}[Y] < k/(1 + \varepsilon)$$

as well (note  $\text{Var}[Y] = \sum_i \text{Var}[Y_i]$  since the  $Y_i$  are pairwise independent). Thus by Chebyshev's inequality,

$$\mathbb{P}(Y \geq k) \leq \mathbb{P}(|Y - \mathbb{E} Y| > k(1 - 1/(1 + \varepsilon))) < \frac{k}{1 + \varepsilon} \cdot \frac{(1 + \varepsilon)^2}{k^2 \varepsilon^2} = \frac{1 + \varepsilon}{\varepsilon^2 k} < 1/6.$$

We can similarly bound the probability that  $\tilde{t} < (1 - \varepsilon)t$ . This can only happen if there weren't *enough* distinct indices in the stream that hashed to small values under  $h$ . Specifically, let  $Z_i$  be an indicator random variable for the event that the  $i$ th distinct integer in the stream hashed to a value below  $kM/((1 - \varepsilon)t)$ , and define  $Z = \sum_i Z_i$ . Then  $\tilde{t} < (1 - \varepsilon)t$  can only occur if  $Z < k$ . But we have  $k/((1 - \varepsilon)t) \geq \mathbb{E} Z_i > k/((1 - \varepsilon)t) - 1/M \geq (1 + \varepsilon)k/t - 1/M$  (the  $1/M$  is due to rounding). We note  $1/M < \varepsilon k/(4t)$  since  $\varepsilon > 1/\sqrt{n}$  and  $t < n$ . Thus  $k/((1 - \varepsilon)t) \geq \mathbb{E} Z_i > (1 + 3\varepsilon/4)k/t$ , implying

$$k/(1 - \varepsilon) \geq \mathbb{E} Z > (1 + 3\varepsilon/4)k$$

and also since  $Z$  is a sum of pairwise independent Bernoullis, again

$$\text{Var}[Z] \leq \mathbb{E} Z \leq k/(1 - \varepsilon).$$

Thus by Chebyshev's inequality,

$$\mathbb{P}(Z < k) < \mathbb{P}(|Z - \mathbb{E} Z| > (3/4)\varepsilon k) < \frac{k}{1 - \varepsilon} \cdot \frac{16}{9\varepsilon^2 k^2} < \frac{16}{9(1 - \varepsilon)} \cdot \frac{1}{\varepsilon^2 k} < 1/6,$$

as desired.

**Other comments.** It is known, via a different algorithm, that for constant failure probability space  $O(1/\varepsilon^2 + \log n)$  is achievable [KNW10b], and furthermore this is optimal [AMS99, Woo04] (also see [TSJ08]). An algorithm that is more commonly used in practice is HyperLogLog [FEFGM07]. Although it assumes access to a truly random hash function, and asymptotically uses a factor  $\lg \lg n$  more space than the algorithm of [KNW10b], its performance in practice is superior. For a historical discussion of the development of HyperLogLog and other distinct elements and approximate counting algorithms (e.g. Morris), see [Lum18]. For arbitrary, potentially subconstant failure probability  $\delta$ , the optimal bound is  $\Theta(\varepsilon^{-2} \log(1/\delta) + \log n)$ ; the upper bound was shown in [Bl20], and the lower bound in [AMS99, JW13].

### 2.2.3 Another algorithm via geometric sampling

We now describe another approach to designing a small-space algorithm for the distinct elements problem. This algorithm will use more memory than the KMV algorithm, but it serves the pedagogical value of introducing a technique that is common in the design of streaming algorithms: *geometric sampling*. The main idea of geometric sampling is to pick a hash function  $h : [n] \rightarrow \{0, \dots, \log n\}$  from a 2-wise independent family with elements in the hash function's range having geometrically decreasing probabilities:  $\mathbb{P}(h(i) = j) = 2^{-(j+1)}$  for  $j < \log n$  (and  $\mathbb{P}(h(i) = \log n) = 1/n$ )<sup>2</sup>. This hash function  $h$  then naturally partitions the stream into  $\log n + 1$  different “substreams”  $S_0, \dots, S_{\log n}$ , where  $S_r$  is the stream restricted to  $\{i : h(i) = r\}$ . Each  $S_r$  is then fed into some separate data structure  $D_r$ . In the specific case of distinct elements,  $D_r$  is a data structure that solves the following promise problem (for some parameter  $k$ ): if  $t \leq k$ , then the number of distinct elements must be output exactly; otherwise, the data structure may fail in an arbitrary way. This problem has a trivial deterministic solution using  $O(k \log n)$  bits of memory: write down the indices of the first  $k$  distinct elements seen (and if a  $(k+1)$ st element is seen, simply write **Error**). We first show that such an algorithm (with  $k = k_1$  a constant) can give a constant-factor approximation to  $t := F_0$  with good probability. We then run the same algorithm in parallel with some larger  $k = k_2 = C/\varepsilon^2$ . Let then  $r^*$  be such that  $t/2^{r^*} = \Theta(1/\varepsilon^2)$ , which we can calculate based on our constant-factor approximation to  $t$  (in the case that  $t \ll 1/\varepsilon^2$  there is no such  $r^*$ , but this case can be solved exactly by keeping track of an extra data structure  $D'$  with  $k = k_2$  that ignores  $h$  and processes all stream updates). In expectation, the number of distinct elements  $F_0(r^*)$  at level  $r^*$  is  $t/2^{r^*}$ , and the variance is upper bounded by this quantity as well. Thus by Chebyshev's inequality,  $F_0(r^*) = t/2^{r^*} \pm O(\sqrt{t/2^{r^*}}) = (1 \pm \varepsilon)t/2^{r^*}$  with large constant probability. By choosing  $C$  large enough, we can also guarantee that  $(1 + \varepsilon)t/2^{r^*} \leq k_2$  so that we typically output  $F_0(2^{r^*})$  correctly, exactly. Then multiplying the output of  $D_{r^*}$  by  $2^{r^*}$  gives  $(1 \pm \varepsilon)t$  with large probability.

**A constant-factor approximation.** The data structure is as mentioned above with  $k_1 = 1$ . Given an integer  $i$  in the stream, we feed it to the data structure  $D_{h(i)}$ . Then to obtain a constant-factor approximation, we return  $2^\ell$  for  $\ell$  being the largest value of  $r \in \{0, \dots, \log n\}$  such that  $D_r.\text{query}() \neq 0$  (even if  $D_r$  reports an error, we count that as not returning 0). Note that  $\ell$  is a random variable, as it depends on  $h$ .

In the analysis below, we define  $L := \log_2 t$ .

**Lemma 2.2.8.**  $\mathbb{P}(\ell > L + 4) \leq 1/8$ .

*Proof.* For any  $r$  we have  $\mathbb{E}_h F_0(r) = t/2^r$ . Thus

$$\begin{aligned} \mathbb{P}(\ell > L + 4) &= \mathbb{P}(\exists r > L + 4 : F_0(r) > 0) \\ &\leq \sum_{r=\lceil L+4 \rceil}^{\log n} \mathbb{P}(F_0(r) > 0) \text{ (union bound)} \\ &= \sum_{r=\lceil L+4 \rceil}^{\log n} \mathbb{P}(F_0(r) \geq 1) \text{ (} F_0(r) \text{ is an integer)} \\ &= \sum_{r=\lceil L+4 \rceil}^{\log n} \mathbb{P}(F_0(r) \geq \frac{2^r}{t} \cdot \mathbb{E} F_0(r)) \end{aligned}$$

<sup>2</sup>One way to implement such  $h$  is to select  $h' : [n] \rightarrow [n]$  from a pairwise independent family and defining  $h(i)$  to be the index of the least significant bit of  $h'(i)$ .

$$\begin{aligned}
&\leq \sum_{r=\lceil L+4 \rceil}^{\log n} \frac{t}{2^r} \\
&< \sum_{q=4}^{\infty} 2^{-q} \\
&= \frac{1}{8}.
\end{aligned}$$

□

**Lemma 2.2.9.**  $\mathbb{P}(F_0(\lceil L - 4 \rceil) = 0) < 1/8$ .

*Proof.* Write  $q = F_0(\lceil L - 4 \rceil)$  so that  $\mathbb{E} q = t/2^{\lceil L-4 \rceil}$ , and  $\text{Var}[q] < \mathbb{E} q$  by pairwise independence of  $h$ . Then

$$\begin{aligned}
\mathbb{P}(q = 0) &\leq \mathbb{P}(|q - \mathbb{E} q| \geq \mathbb{E} q) \\
&< \mathbb{P}(|q - \mathbb{E} q| \geq \sqrt{\mathbb{E} q} \sqrt{\text{Var}[q]}) \\
&\leq \frac{1}{\mathbb{E} q} \text{ (Chebyshev's inequality),}
\end{aligned}$$

which is less than  $1/8$  by the definition of  $L$ . □

**Lemma 2.2.8** tells us that it is unlikely that  $\ell > \log_2 + 4$  (and thus unlikely that our estimator is bigger than  $16t$ ). **Lemma 2.2.9** implies it is unlikely that  $\ell < \lceil L - 4 \rceil$ . Thus by a union bound, we have the following result.

**Lemma 2.2.10.**  $2^\ell \in [\frac{1}{16}t, 16t]$  with probability at least  $3/4$ .

**Refining to  $(1 + \varepsilon)$ -approximation.** We run the constant-factor approximation algorithm in parallel, so that at any point we can query it to obtain  $\hat{t} := 2^\ell$ . We henceforth condition on the good event that  $\hat{t} \in [\frac{1}{16}t, 16t]$ .

Define  $k_2 = \lceil 512/\varepsilon^2 + 68/\varepsilon \rceil$ . We first check whether  $t \leq k_2$  by checking  $D'$ ; if so, we return  $D'.\text{query}()$ , which is guaranteed to equal  $t$  exactly with probability 1. Otherwise, as mentioned in **Subsection 2.2.3** we define  $r^*$  such that  $t/2^{r^*} = \Theta(1/\varepsilon^2)$ ; specifically we set  $r^* = \lceil \log_2(\varepsilon^2 \hat{t}/32) \rceil$ . Then due to rounding from the ceiling function and also by our guarantees on  $\hat{t}$ , we have  $1/\varepsilon^2 < t/2^{r^*} \leq 512/\varepsilon^2$ . We also have  $\text{Var}_h[F_0(r^*)] \leq t/2^{r^*}$ , and thus by Chebyshev's inequality  $F_0(r^*) = t/2^{r^*} \pm 3\sqrt{512}/\varepsilon = t/2^{r^*} \pm 68/\varepsilon = (1 \pm 68\varepsilon)t/2^{r^*}$  with probability at least  $8/9$ . In this case, (1) scaling up  $F_0(r^*)$  by  $2^{r^*}$  gives a  $1 + O(\varepsilon)$  approximation as desired (which can be made  $1 + \varepsilon$  by running this algorithm with  $\varepsilon' = \varepsilon/68$ , and (2) we can know  $F_0(r^*)$  exactly (in order to scale it up) since  $F_0(r^*) \leq k_2$  by our choice of  $k_2$ . In summary we have:

**Theorem 2.2.11.** *There is an algorithm for  $(1 + \varepsilon)$ -approximation of  $F_0$  with probability  $1 - \delta$  using  $O(\varepsilon^{-2} \log^2 n \log(1/\delta))$  bits of memory.*

*Proof.* The constant factor approximation succeeds with probability at least  $3/4$ . Conditioned on it succeeding, our refinement to  $1 + \varepsilon$  approximation succeeds with probability at least  $8/9$ . Thus our overall success probability is at least  $3/4 \cdot 8/9 = 2/3$ , which can be amplified to  $1 - \delta$  in the usual manner by returning the median of  $\Theta(\log(1/\delta))$  copies of this basic data structure.

We now analyze the memory to obtain success probability  $2/3$ ; for probability  $1 - \delta$ , this increases by a multiplicative  $\log(1/\delta)$  factor. Recall we run two algorithms in parallel, with  $k_1 = 1$

and  $k_2 = \Theta(1/\varepsilon)^2$ ; we focus on the latter since its memory consumption is strictly more. Storing the hash function  $h$  requires only  $O(\log n)$  bits.  $D'$  takes an additional  $O(k_2 \log n) = O(\varepsilon^{-2} \log n)$  bits of memory. Similarly, each  $D_i$  takes  $O(k_2 \log n)$  bits of memory and there are  $\log_2 n + 1$  such structure, yielding the desired space bound in total.  $\square$

**Remark 2.2.12.** One advantage of this algorithm over the KMV algorithm is that its update and query times are faster; in fact in the so-called “transdichotomous word RAM model” (where the word size is at least  $\log n$  bits) [FW93], these times are  $O(1)$  (for success probability  $2/3$ ), whereas the natural implementation of KMV would use a heap and have  $\Theta(\log(1/\varepsilon))$  update time. To see this, we can store the distinct elements in  $D'$  and each  $D_r$  using a dynamic dictionary data structure such as hashing with chaining, which supports  $O(1)$  worst case expected update time. We also know which  $D_r$  to update since we can compute  $h(i)$  for each  $i$  in constant time. For the constant-factor approximation algorithm, since  $k_1 = 1$  we need not store the identity of the seen item in  $D_r$ ; we can simply mark a bit 0 or 1 as to whether any item has ever been seen. As there are  $\log n$  such bits, we can store them in a single machine word  $x$ . We can find the largest non-empty  $D_r$  by computing the index of the least significant set bit of  $x$ , which can be done in  $O(1)$  time in the word RAM model [Bro93]. For the parallel data structure with larger  $k_2 = \Theta(1/\varepsilon^2)$ , updates are similarly fast; for queries, we only query the size of a single  $D_r$  (namely  $D_{r^*}$ ), which takes constant time. Thus overall, queries are worst case  $O(1)$  time, and updates are  $O(1)$  expected time (the bottleneck being the use of dynamic dictionary in each  $D_r$ ).

## 2.3 Quantiles

In this problem we see a stream of comparable items  $y_1, y_2, \dots, y_n$ , in some arbitrary order. For example,  $x_i$  could be the response time to serve request  $i$  in some system. We will use  $x_1, \dots, x_n$  throughout this section to denote the  $y_i$  in sorted order:  $x_1 < x_2 < \dots < x_n$ . Without loss of generality we can assume no two items are equal, since any algorithm can interpret  $y_i$  in the stream as  $(y_i, i)$  and use lexicographic ordering. We are also given up front some  $\varepsilon \in [0, 1)$  at the beginning of the stream. At query time we would then like to support three (really two) types of queries:

- **rank( $x$ ):** if  $r$  is such that  $x_r < x \leq x_{r+1}$  (i.e. it is the true rank of  $r$ ), then return  $r \pm \varepsilon n$ .
- **select( $r$ ):** return any item whose rank in the sorted order is  $r \pm \varepsilon n$ .
- **quantile( $\phi$ ):** this is syntactic sugar for **select( $\phi n$ )**.

Going back to the example, a typical use of such a data structure is to answer quantile queries for  $\phi = 0.999$  say, when monitoring the performance of a system to ensure that 99.9% of user queries are responded to quickly. We remark that the assumption that the  $y_i$  are distinct is not necessary but just simplifies the remaining discussion in this section; without that assumption, one could leave the requirement of **rank( $x$ )** unchanged (any  $r$  satisfying the given inequalities is acceptable), and for **select( $r$ )**, if  $x_a = x_{a+1} = \dots = x_b = x$ , then it is acceptable to return  $x$  as long as  $[a, b] \cap [r - \varepsilon n, r + \varepsilon n] \neq \emptyset$ .

The state-of-the-art algorithms for this problem are as follows. By a “comparison-based” sketch, we mean a sketch that works in the model where items come from an infinite universe in which between any two values  $x < y$  there exists a  $z$  such that  $x < z < y$ , and the memory of the algorithm can be divided into  $M_1$  and  $M_2$  (see [CV20]). Here  $M_1$  is a subset of items seen, and  $M_2$  is some extra storage that is only allowed to keep track of results of comparisons between stream items and those items in  $M_1$ . When processing a new stream update  $x$ , we are allowed to compare

$x$  with all elements of  $M_1$ , possibly store some of those results in  $M_2$ , then also possibly include  $x$  in  $M_1$  and remove some current items from  $M_1$ . All decisions made are determined only by the results of comparisons.

- **Deterministic, not comparison-based.** If the input values are known to come from a finite universe  $\{0, \dots, U-1\}$ , the **q-digest** sketch [SBAS04] solves quantiles using  $O(\varepsilon^{-1} \log U)$  words of memory.
- **Deterministic, comparison-based.** The GK sketch [GK01] uses  $O(\varepsilon^{-1} \log(\varepsilon n))$  words of memory, which was recently shown to be best possible for any deterministic comparison-based algorithm [CV20].
- **Randomized, comparison-based.** The KLL sketch [KLL16] when given parameter  $\delta$  at the beginning of the stream uses  $O(\varepsilon^{-1} \log \log(1/\delta))$  words of memory and succeeds with probability  $1 - \delta$ . That is

$$\forall \text{ queries } q, \mathbb{P}(\text{KLL answers } q \text{ correctly}) \geq 1 - \delta$$

If one wants the “all-quantiles” guarantee, i.e.

$$\mathbb{P}(\forall \text{ queries } q, \text{ KLL answers } q \text{ correctly}) \geq 1 - \delta,$$

then this is accomplished using  $O(\varepsilon^{-1} \log \log(1/(\varepsilon \delta)))$  words of memory. For achieving the first guarantee, they also prove that their bound is optimal up to a constant factor (see also [CV20] for a strengthening of the lower bound statement).

**Open:** Is **q-digest** optimal?

**Open:** Can one do better than the KLL sketch by not being comparison-based?

**Open:** The GK analysis is complex; can one match its performance via a simpler approach?

For more details of these and several other sketches, see the survey by Greenwald and Khanna [GK16]. As that survey was written in 2016 though, some of the more recent results on quantiles are not covered, e.g. the KLL sketch [KLL16] and a faster version using the same memory [ILL<sup>+</sup>19], an optimal lower bound for deterministic comparison-based sketches [CV20], and some recent works on relative-error quantile sketches (see [CKL<sup>+</sup>20] and the references therein). Some of these more recent developments appear in Chapter 4 of the upcoming book of Cormode and Yi [CY20].

In the following sections, we will cover the following sketches: (1) the **q-digest** sketch, (2) the MRL sketch [MRL98], which is deterministic and comparison-based but achieves the worse memory bound  $O(\varepsilon^{-1} \log^2(\varepsilon n))$  when compared to the GK sketch, and (3) KLL.

### 2.3.1 q-digest

In this setting, where items have values in  $\{0, \dots, U-1\}$ , we will allow for items of the stream to have the same value and not do the reduction  $y_i \mapsto (y_i, i)$  mentioned in Section 2.3. The main idea of this data structure is to conceptually imagine a perfect binary tree whose leaves correspond to the values  $0, 1, \dots, U-1$ . For every node  $v$  in this tree, we have a counter  $c[v]$ . Each node represents an interval of values, from its leftmost leaf descendant to its rightmost. The **q-digest** structure will store the names of all  $v$  with  $c[v] \neq 0$ , together with their  $c[v]$  values. Ideally we would like that  $c[v]$  is zero for all internal nodes of the tree, and equals the number of occurrences of  $v$  in the stream for every leaf node  $v$ . Doing so is equivalent to simply storing a histogram: for each  $x \in \{0, \dots, U-1\}$ , we keep track of the number of stream elements equal to  $x$ . An exact histogram



would of course allow us to answer all rank/select queries exactly (with  $\varepsilon = 0$ ), but unfortunately could be too memory-intensive if too many distinct values are seen in the stream (which could be as large as  $\min\{n, U\}$  values). To remedy this, **q-digest** occasionally merges two sibling nodes into their parent, by zeroing out their  $c[]$  values after adding them to their common parent. This naturally saves memory by zeroing out nodes, at the expense of worsening accuracy: for counts stored at an internal nodes  $v$ , we cannot be sure what precise values the items contributing to those counts had within  $v$ 's subtree.

In what follows we assume  $U$  is a power of 2, and we write  $L := \log_2 U$ . This is without loss of generality, as if not we could simply round  $U$  upward and increase it by at most a factor of 2. In the perfect binary tree, we refer to the sibling of a non-root node  $v$  as  $s[v]$ , and its parent as  $p[v]$ . For a non-leaf  $v$ , we use  $\text{left}[v]$ ,  $\text{right}[v]$  to denote its left and right children, respectively.

```

q-digest sketch:
initialization //  $\varepsilon \in (0, 1)$  is given
1.  $n \leftarrow 0$ 
2.  $S \leftarrow \emptyset$  // set of  $(v, c[v])$  pairs with  $c[v] \neq 0$ ; note  $c[v] = 0$  for  $v$  not tracked by  $S$ 
3.  $K := \lceil 6L/\varepsilon \rceil$ 

insert( $i$ ) // see item  $i$  in the stream
1.  $n \leftarrow n + 1$ 
2. if leaf node  $i$  is in  $S$ :
3.    $c[i] \leftarrow c[i] + 1$ 
4. else:
5.    $S \leftarrow S \cup \{(i, 1)\}$ 
6.   if  $|S| > K$ :
7.     compress()

compress() // the root has depth 0, and leaves are at depth  $L$ 
1. for  $\ell = L, L - 1, \dots, 1$ :
2.   for each node  $v$  in  $S$  at level  $\ell$ :
3.     if  $c[p[v]] + c[v] + c[s[v]] \leq \varepsilon n/L$ :
4.        $c[p[v]] \leftarrow c[p[v]] + c[v] + c[s[v]]$ 
5.        $c[v] \leftarrow 0$ 
6.        $c[s[v]] \leftarrow 0$ 
7.       adjust  $S$  as needed (remove  $v$  and  $s[v]$ , and associate  $p[v]$  with its new  $c[p[v]]$  value)

rank( $x$ ) // let  $t[v]$  denote the the sum of all  $c[u]$  values of nodes  $u$  in  $v$ 's subtree
1. // below we express  $x = \sum_{i=0}^{L-1} x_i 2^i$  in binary
2.  $v \leftarrow \text{root}$ 
3.  $\text{ans} \leftarrow c[v]$ 
4. for  $i = L - 1, L - 2, \dots, 0$ :
5.   if  $x_i = 1$ :
6.      $\text{ans} \leftarrow \text{ans} + t[\text{left}[v]]$ 
7.      $v \leftarrow \text{right}[v]$ 
8.   else:
9.      $v \leftarrow \text{left}[v]$ 
10.   $\text{ans} \leftarrow \text{ans} + c[v]$ 
11. return  $\text{ans}$ 

select( $r$ )
1.  $A \leftarrow c[\text{root}]$ ,  $v \leftarrow \text{root}$ ,  $x \leftarrow 0$ 
2. for  $i = L - 1, L - 2, \dots, 0$ :
3.   if  $A + t[\text{left}[v]] \geq r - 1$ :
4.      $v \leftarrow \text{left}[v]$ 
5.   else:
6.      $A \leftarrow A + c[\text{left}[v]]$ 
7.      $v \leftarrow \text{right}[v]$ 
8.      $x \leftarrow x + 2^i$ 
9.    $A \leftarrow A + c[v]$ 
10. return  $x$ 

```

Figure 2.1: Pseudocode for **q-digest**; bottom-up implementation.

Fig. 2.1 provides a bottom-up implementation of `q-digest` (bottom-up since updates directly affect the leaves, and only percolate upward later when `compress` is called).  $S$  can be implemented as any dynamic dictionary data structure. A top-down recursive implementation of both `insert` and `compress` is also possible; see [CY20, Chapter 4] for details.

To answer queries, we let  $\tau[v]$  denote the sum of all  $c[\cdot]$  values of nodes in its subtree, including its own  $c[v]$  value. When answering `rank(x)`, we return the sum of all  $c[\cdot]$  values of leaves to the left of and including  $x$ , as well as their ancestors. We can accomplish this by traversing the unique root to  $x$  path in the tree and keeping a running sum (see Fig. 2.1). To answer `select(r)`, we would like to find the unique value  $x \in \{0, \dots, U - 1\}$  such that the sum of all  $c[\cdot]$  values of leaves to the left of and including  $x$ , and their ancestors, is at least  $r - 1$ , whereas the same computation for  $x - 1$  would be strictly less than  $r - 1$ . We find such  $x$  via a depth-first search (DFS) from the root. We keep a running sum  $A$  that approximates the number of stream elements that are strictly less than the entire range covered by the current node  $v$ , and we return the leaf that our DFS lands; see Fig. 2.1 for details.

**Lemma 2.3.1.** *`q-digest` answers both `rank` and `select` correctly, i.e. with error  $\pm \epsilon n$  for each.*

*Proof.* The key observation is that any element  $x$  in the stream contributes to the  $c[\cdot]$  value of exactly one node: either  $x$ 's leaf node, or one of its ancestors. For the calculation of `rank(x)`, note that any  $z \leq x$  either increments the counter of some ancestor of  $x$ , or of some node in the left subtree of an ancestor of  $x$ , and thus it is counted. For  $z > x$ , it either increments the counter of some node in the right subtree of an ancestor of  $x$  (in which case it is not counted when we calculate `rank(x)`), or it increments the counter of an ancestor of  $x$  (in which case it is counted). But  $x$  has only  $L$  ancestors, and each one has a counter value of at most  $\epsilon n/L$  by `compress`, so the total error introduced by larger  $z$  is at most  $L \cdot (\epsilon n/L) = \epsilon n$ . Similar reasoning implies the correctness of `select`, which we leave as an exercise to the reader.  $\square$

The following lemma motivates our choice of  $K$  in the pseudocode, as it means we only have to call `compress` after every at least roughly  $K/2$  updates. Thus the amortized runtime to `compress` is improved compared with running it after every insertion (which would still be correct and improve the space by a factor of two).

**Lemma 2.3.2.** *After any call to `compress`, the resulting  $S$  will have size at most  $3L/\epsilon + 1$ .*

*Proof.* Consider  $S$  after compression. Then for any node  $v$  tracked by  $S$  which is not the root,

$$c[p[v]] + c[v] + c[s[v]] > \epsilon n/L.$$

Summing over such  $v$ ,

$$3n \geq \sum_{\substack{v \in S \\ v \text{ not the root}}} c[p[v]] + c[v] + c[s[v]] > (|S| - 1)\epsilon n/L, \quad (2.9)$$

where the LHS holds since each stream element contributes to exactly one  $c[v]$  counter, and each  $c[v]$  appears at most three times in the above sum (it is the sibling of at most one node in  $S$ , and the parent of at most one node in  $S$ ). Rearranging gives  $|S| \leq 3L/\epsilon + 1$ , as desired.  $\square$

**Remark 2.3.3.** There are several valid decision choices in an implementation of `q-digest`. For example, it is correct to call `compress` after every insertion, though that would increase runtime unnecessarily. Alternatively, one could percolate up merges into the parent starting from node  $i$  until reaching either the root or a parent whose  $c[\cdot]$  value would be too large to allow the merge.

Doing so helps deamortize the data structure to prevent frequent calls to `compress` (every  $\approx K/2$  updates), but at the same time, note that  $n$  is changing: it increases after every insertion. Thus capacities of nodes to hold more in their counters is continuously increasing: previous parents which rejected merges may later allow them due to increased capacity. To ensure that  $|S|$  always stays small, one could then call `compress` every time  $n$  doubles, say. The RHS of Eq. (2.9) in between calls to `compress` would still be at least  $(|S| - 1)\varepsilon n / (2L)$ , implying that the data structure always has size at most  $6L/\varepsilon + 1$  even between calls to `compress`.

### 2.3.2 MRL

**MRL sketch:**

**initialization** //  $\varepsilon \in (0, 1)$  and the final stream length  $n$  are given

1.  $k := \varepsilon^{-1} \lceil \log(\varepsilon n) \rceil + 1$ , rounded up to the nearest even integer
2.  $L := \lceil \log(n/k) \rceil$
3. // the below  $A_j$  arrays are 1-indexed
4. initialize empty “compactor” arrays  $A_0, \dots, A_L$  each of size  $k$ , initialized to have all null entries

**insert( $i$ )** // see item  $i$  in the stream

1. insert  $i$  into  $A_0$  // i.e. replace the first non-null entry in  $A_0$  with  $i$
2.  $j \leftarrow 0$
3. **while**  $A_j$  has size  $k$ : // i.e. no non-null entries
4.      $\text{sort}(A_j)$
5.     insert  $A_j[1], A_j[3], \dots, A_j[n-1]$  into  $A_{j+1}$
6.     set all entries of  $A_j$  to null
7.      $j \leftarrow j + 1$

**rank( $x$ )**

1.  $\text{ans} \leftarrow 0$
2. **for**  $j = 0, 1, \dots, L$ :
3.     **for each** item  $z$  in  $A_j$ :
4.         **if**  $z < x$ :
5.              $\text{ans} \leftarrow \text{ans} + 2^j$
6. **return ans**

**select( $r$ )**

1.  $B \leftarrow$  array containing  $(z, 2^j)$  for every  $z \in A_j$ , over all  $j \in \{0, \dots, L\}$
2. sort  $B$  lexicographically
3.  $i \leftarrow$  smallest index such that  $\sum_{j=1}^i B[j].\text{second} \geq r$
4. **return**  $B[i].\text{first}$

Figure 2.2: Pseudocode for MRL sketch.

The MRL data structure of [MRL98] is also deterministic, but is comparison-based and therefore does not need to assume knowledge of where the universe comes from. It does however need to know the length  $n$  of the final stream in advance<sup>3</sup>. The space complexity will be  $O(\varepsilon^{-1} \log^2(\varepsilon n))$ . Though we have not yet spoken about *mergeability* of sketches, we briefly do so now to state an open problem related to this concept. Roughly speaking, a sketching algorithm is *fully mergeable* if given two sketches  $S_1$  and  $S_2$  created from inputs  $D_1$  and  $D_2$ , a sketch  $S$  of  $D := D_1 \sqcup D_2$  can be created with no degradation in quality of error or failure probability, and satisfying the same efficiency constraints as  $S_1, S_2$ . Furthermore, this mergeability condition should hold in an arbitrary merge tree of several sketches (i.e. it should be possible to merge a previous merger of sketches with another merger of sketches, etc.). The MRL sketch unfortunately does not satisfy this property because, as we will soon see, the data structure sets an internal parameter  $k$  based on  $n$  at the beginning of the sketching procedure. Thus if  $k_1$  is set based on  $n_1 := |D_1|$  and  $k_2$  based

<sup>3</sup>The more space-efficient GK sketch does not need to know  $n$  in advance. Also, for the MRL sketch it suffices to only know an upper bound  $N$  on the final stream length, and the memory bounds will then have  $n$  replaced with  $N$ .

on  $n_2 := |D_2|$ , it is not clear how to combine the sketches with different  $k$  parameters into a new sketch with  $k_3 = n_1 + n_2$  without increasing the error  $\varepsilon$  guaranteed by the original two sketches.

**Open:** Does there exist a deterministic, comparison-based quantiles sketch using  $o(n)$  memory which is fully mergeable? It is conjectured in [ACH<sup>+</sup>13] that the answer is no.

The main component in the MRL sketch is a *compactor*, which has a single parameter  $k$  that is a positive even integer. A compactor stores anywhere between 0 and  $k$  items. When it is full, i.e. has  $k$  items, it “compacts”: this amounts to sorting its elements into  $x_1 < x_2 < \dots < x_k$  then outputting the odd-indexed elements  $x_1, x_3, x_5, \dots, x_{k-1}$ . In the MRL sketch compactors are chained together, so that the output of one compactor is inserted into the next compactor, possibly causing a chain reaction; see Fig. 2.2.

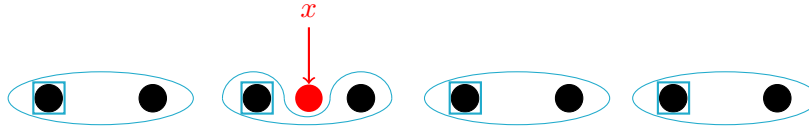


Figure 2.3: Illustrative figure for understanding error introduced during compaction.

A picture to keep in mind for how compaction introduces error into *rank* queries is drawn in Fig. 2.3. Imagine the query element is  $x$  (the red dot), and let the black dots be the elements of the stream currently living in one of the compactors with  $k = 8$ , arranged in sorted order. In truth, 3 of these elements are less than  $x$ . But after compaction and moving the four boxed elements into the next level, we will see that two of the boxed elements are less than  $x$  and thus think that 4 original elements were less than  $x$ . Being off by one in this count introduces an error of  $2^j$  for compactions performed at level  $j$ . Note also that if  $x$  had even rank in the sorted order amongst these elements, no error would be introduced at all. Thus the error introduced for any  $x$  during a compaction at level  $j$  can be expressed as  $\eta_x 2^j$ , where  $\eta_x \in \{0, 1\}$  depends on the parity of  $x$ 's rank amongst the items that were compacted. Thus for any  $x$ , querying the rank of  $x$  has additive error

$$E(x) := \sum_{j=0}^{L-1} \sum_{i=1}^{m_j} \eta_{x,j,i} 2^j, \quad (2.10)$$

where  $m_j$  is the number of compactions performed by  $A_j$ . Since there are only  $n$  items in the stream, there are at most  $n/2^j$  items that are ever inserted into  $A_j$ . A compaction clears out  $k$  elements at a time, and thus  $m_j \leq n/(k2^j)$ . Therefore

$$\begin{aligned} E(x) &\leq \sum_{j=0}^{L-1} \sum_{i=1}^{m_j} 2^j \\ &\leq \sum_{j=0}^{L-1} \frac{n}{k2^j} 2^j \\ &= \frac{n}{k} \cdot L \end{aligned} \quad (2.11)$$

To ensure Eq. (2.11) is at most  $\varepsilon n$ , we must set  $k$  so that  $L/k = \lceil \log(n/k) \rceil / k \leq \varepsilon$ . That is,  $k \geq \varepsilon^{-1} \lceil \log(n/k) \rceil$ . This inequality is satisfied by our choice of  $k$ .

**Theorem 2.3.4.** *The MRL sketch uses space  $O(k \log(n/k)) = O(\varepsilon^{-1} \log^2(\varepsilon n))$ .*

### 2.3.3 KLL

KLL is a randomized sketch that, for any fixed rank query, succeeds with probability  $1 - \delta$ . Its space usage is  $O(\varepsilon^{-1} \log \log(1/\delta))$  words. As mentioned in [Section 2.3](#), it can also provide the “all-quantiles” guarantee using  $O(\varepsilon^{-1} \log \log(1/(\varepsilon\delta)))$ , in which case `select` can be implemented by performing rank queries on every stored item.

The starting point of KLL is a randomized improvement to the MRL sketch given in [\[ACH<sup>+</sup>13\]](#). This randomized MRL sketch is identical to the MRL sketch, except that rather than a compactor always output the odd-indexed elements, it uniformly at random decides to either output the odd- or even-indexed elements in the sorted order. Thus, again looking at [Fig. 2.3](#), we see that the error introduced when estimating  $\text{rank}(\mathbf{x})$  by any compaction at level  $j$  can be expressed as  $\eta\sigma 2^j$ , where  $\eta \in \{0, 1\}$ , and  $\sigma \in \{-1, 1\}$  is uniformly random. Then we can rewrite the error term in our estimate for  $\text{rank}(\mathbf{x})$  as

$$E(\mathbf{x}) := \sum_{j=0}^{L-1} \sum_{i=1}^{m_j} \eta_{\mathbf{x},j,i} \sigma_{j,i} 2^j, \quad (2.12)$$

This error random variable precisely fits the setup of Khintchine’s inequality (see [Theorem 1.1.7](#)), where the vector “ $x$ ” in the statement of the inequality has entries  $(\eta_{\mathbf{x},j,i} 2^j)_{j,i}$ . Then we see that

$$\|x\|_2^2 \leq \sum_{j=0}^{L-1} m_j 2^{2j} \leq \frac{n}{k} \cdot 2^L \leq 2 \left(\frac{n}{k}\right)^2. \quad (2.13)$$

Plugging into Khintchine’s inequality, we therefore have

$$\mathbb{P}(|E(\mathbf{x})| > \varepsilon n) \leq 2e^{-\varepsilon^2 n^2 \cdot \frac{k^2}{2n^2}}, \quad (2.14)$$

which is at most  $\delta$  for  $k = \lceil \varepsilon^{-1} \sqrt{2 \log(1/\delta)} \rceil$ . Note also a further benefit: since  $k$  does not depend on  $n$ , we do not need to know  $n$  in advance. Although the number of compactors does depend on  $n$ , we can simply start off with only one compactor and allocate new compactors as they become needed. We thus have the following theorem.

**Theorem 2.3.5.** *There is a randomized, comparison-based sketch for quantiles with additive error  $\pm \varepsilon n$  and failure probability at most  $\delta$  using  $O(\varepsilon^{-1} \sqrt{\log(1/\delta)} \log(\varepsilon n / \sqrt{\log(1/\delta)}))$  words of memory. This algorithm does not need to know the stream length  $n$  in advance.*

KLL then makes further optimizations to improve the space complexity. The first optimization is to not make all the compactors have equal size. Rather, the last compactor  $\mathbf{A}_L$  will have (the largest) size  $k$ . Then compactor  $\mathbf{A}_{L-j}$  will have size equal to the maximum of 2 and  $\approx (2/3)^j k$ . Since we do not know  $n$  a priori and thus do not know how many compactors we will need, this is achieved by making our compactors have dynamically changing sizes. For example, initially we will only have one compactor  $\mathbf{A}_0$  with capacity  $k$ . Then when a compaction finally happens, we will create  $\mathbf{A}_1$  with capacity  $k$  and shrink the capacity of  $\mathbf{A}_0$  to its new, smaller capacity. It can be shown that the number of compactors at any given time is still  $\log(n/k) + O(1)$ . Furthermore, the number of compactions  $m_j$  at level  $j$  is still at most  $n/(k2^j)$ . Similar computations to [Eqs. \(2.13\) and \(2.14\)](#) then imply  $k = O(\varepsilon^{-1} \sqrt{\log(1/\delta)})$  still suffices, but since our compactor sizes are geometrically decreasing (down to a minimum size of 2), our space is improved. In particular, we have the following theorem, which improves [Theorem 2.3.5](#) by making the  $\log n$  term additive instead of multiplicative.

**Theorem 2.3.6.** *There is a randomized, comparison-based sketch for quantiles with additive error  $\pm \varepsilon n$  and failure probability at most  $\delta$  using  $O(\varepsilon^{-1} \sqrt{\log(1/\delta)} + \log(\varepsilon n / \sqrt{\log(1/\delta)}))$  words of memory. This algorithm does not need to know the stream length  $n$  in advance.*

We only sketch the remaining improvements to the KLL sketch. The first is to observe that since we enforce that all capacities must be at least 2 and otherwise compactor capacities decay geometrically, only the top  $O(\log k)$  compactors have size more than 2, and the bottom  $T = L - O(\log k)$  compactors all have size exactly 2. One can view the action of these bottom compactors then in unison: from the input stream, elements come in, and the output then of compactor  $A_{T-1}$  is a uniformly random element amongst  $2^T$  input stream elements. Thus, rather than spend  $O(T)$  space implementing these  $T$  elements, we can implement a sampler using  $O(1)$  words of memory. This improves the overall space to  $O(k)$  down from  $O(k + \log(n/k))$ . The next idea is to make all the top  $s$  compactors have size  $k$ , and only compactors  $A_{L-j}$  for  $j > s$  have size  $\approx (2/3)^j k$ . Then the space is  $O(sk)$ , but re-examining the error term:

$$\begin{aligned} |E(\mathbf{x})| &= \left| \sum_{j=0}^{L-1} \sum_{i=1}^{m_j} \eta_{\mathbf{x},j,i} \sigma_{j,i} 2^j \right| \\ &\leq \left| \sum_{j=0}^{L-s} \sum_{i=1}^{m_j} \eta_{\mathbf{x},j,i} \sigma_{j,i} 2^j \right| + \left| \sum_{j=L-s+1}^{L-1} m_j 2^j \right| \\ &\leq \underbrace{\left| \sum_{j=0}^{L-s} \sum_{i=1}^{m_j} \eta_{\mathbf{x},j,i} \sigma_{j,i} 2^j \right|}_{\alpha} + \underbrace{\frac{ns}{k}}_{\beta} \end{aligned}$$

Then for the  $\alpha$  term, a simpler computation to Eq. (2.13) yields a bound on  $\|x\|_2^2$  of  $2(n/k)^2 \cdot 2^{-s}$ . We can then apply Khintchine as above to obtain failure probability

$$\mathbb{P}(|\alpha| > \varepsilon n/2) \leq 2e^{-\varepsilon^2 n^2 \cdot \frac{k^2}{8n^2} \cdot 2^s}.$$

If  $s = \lceil \ln \ln(1/\delta) \rceil$ , then the above is at most  $\delta$  as long as  $k \geq \sqrt{8}/\varepsilon$ . On the other hand, we must also ensure  $\beta \leq \varepsilon n/2$ , which requires  $k > s/\varepsilon$ . Thus overall, we can afford to set  $k = O(\varepsilon^{-1} \log \log(1/\delta))$ , which leads to a space bound of  $ks = O(\varepsilon^{-1} (\log \log(1/\delta))^2)$ .

The final optimization to achieve the improved space bound is to observe that the bottleneck above comes from the top  $s$  compactors. We can instead replace them by an optimal deterministic sketch, such as the GK sketch. Thus outputs from  $A_{L-s}$  are fed into the GK sketch and can be viewed as all having the same weight  $2^{L-s+1}$ . When we add in the GK rank query output into our rank query outputs, we thus multiply its output by  $2^{L-s+1}$  first. We thus overall have the following theorem.

**Theorem 2.3.7.** *There is a randomized, comparison-based sketch for quantiles with additive error  $\pm \varepsilon n$  and failure probability at most  $\delta$  using  $O(\varepsilon^{-1} \log \log(1/\delta))$  words of memory. This algorithm does not need to know the stream length  $n$  in advance.*

**Remark 2.3.8.** Although Theorem 2.3.7 gives an optimal memory bound for this problem amongst comparison-based solutions, it is not known to be fully mergeable (see the definition sketch in Subsection 2.3.2). This is because the GK sketch it uses in the last optimization is not known to be fully mergeable. Without this optimization though, KLL is fully mergeable. Thus, the best fully mergeable sketch we have uses space  $O(\varepsilon^{-1} (\log \log(1/\delta))^2)$ ; see [KLL16] for details.

# Chapter 3

## Lower Bounds

In this chapter we discuss common techniques for proving lower bounds for the size of a sketch to solve some problem, or for the memory used by a streaming algorithm to solve some task. There are two main techniques common in this area: compression-based arguments, and communication complexity.<sup>1</sup> We discuss these both in the below sections, together with example applications.

### 3.1 Compression-based arguments

#### 3.1.1 Distinct elements

Recall in [Section 2.2](#) we discussed computing the number of distinct elements in a data stream. In that problem we saw a stream  $x_1, \dots, x_m$  of (possibly non-distinct) elements in a stream, each in  $\{1, \dots, n\}$ . We developed randomized, approximate algorithm: that is, if the true number of distinct elements is  $t$ , the algorithms we discussed achieving non-trivially low memory used randomness to output a number  $\tilde{t}$  satisfying

$$\mathbb{P}(|t - \tilde{t}| > \epsilon t) < \delta.$$

for some  $\epsilon, \delta \in (0, 1)$ . One could ask: were *both* randomization and approximation necessary? Could we have a randomized exact algorithm that fails with some small probability  $\delta$ ? Or a deterministic approximate algorithm? It turns out that the answer to both of these questions is no [[AMS99](#)]. Before proving so, we first show that the strongest possible guarantee, i.e. an algorithm that is both exact *and* deterministic, requires  $\Omega(n)$  bits of memory. The proof we give will be the first example of a *compression argument*. Applied to streaming algorithms, essentially such proofs fit the following template: if there exists an algorithm that uses very little memory (say  $S$  bits), then there exists an injection (we often refer to as the *encoding*)  $\text{Enc} : A \rightarrow \{0, 1\}^{g(S)}$ . Therefore  $g(S) \geq \log_2 A$ , which depending on  $g$ , may imply some lower bound on  $S$  in terms of  $\log A$ . Such arguments are called compression arguments because they show that a streaming algorithm existed that was too good to be true (uses too little memory), then we would be able to compress a big set, i.e.  $A$ , into a smaller one.

**Theorem 3.1.1.** *Suppose  $\mathcal{A}$  is a deterministic streaming algorithm that computes the number of distinct elements exactly. Then  $\mathcal{A}$  uses at least  $n$  bits of memory.*

*Proof.* We show that the existence of  $\mathcal{A}$  implies the existence of an injection  $\text{Enc} : \{0, 1\}^n \rightarrow \{0, 1\}^S$ . Therefore  $S \geq n$ . The injection is defined as follows. Given some  $x \in \{0, 1\}^n$ , we artificially create

---

<sup>1</sup>These techniques are not actually disjoint; some communication complexity bounds are ultimately proven via compression-based arguments.

a stream that contains in some fixed (say sorted) order all  $i$  such that  $x_i = 1$ . We then run the algorithm  $\mathcal{A}$  on it, then define  $\text{Enc}(x)$  to be the memory contents of  $\mathcal{A}$ .

To show that  $\text{Enc}(x)$  is an injection, we show that we can invert it (or “Decode”) via some inverse function  $\text{Dec}$ . If  $M$  is an  $S$  bit string, then  $\text{Dec}(M)$  is defined as follows. We initialize  $\mathcal{A}$  with  $M$  as its memory. We then execute the following:

```

s ←  $\mathcal{A}.\text{query}()$  // support size of  $x$ , i.e.  $|\{i : x_i \neq 0\}|$ 
 $x \leftarrow (0, 0, \dots, 0)$ 
for  $i = 1 \dots n$ :
   $\mathcal{A}.\text{update}(i)$  // append  $i$  to the stream
   $r \leftarrow \mathcal{A}.\text{query}()$  // will either be  $s$  or  $s + 1$ 
  if  $r = s$ : // Encoder must have included  $i$ , so it wasn't a new distinct element
     $x_i \leftarrow 1$ 
   $s \leftarrow r$ 
return  $x$ 

```

□

We next show that no deterministic exact algorithm exists using  $o(n)$  bits of memory. Before doing so though, we introduce the concept of an error-correcting code.

**Definition 3.1.2.** An *error-correcting code* is a set  $\mathcal{C} \subset [q]^t$ . Here  $q$  is referred to as the *alphabet size* and  $\ell$  as the *block length*. The *Hamming distance* between two elements of  $[q]^t$  is the number of entries where they differ, i.e.  $\Delta(x, y) := |\{i : x_i \neq y_i\}|$ . The *relative Hamming distance* is  $\delta(x, y) = \Delta(x, y)/\ell$ . The *distance* of the code is  $\min_{c, c' \in \mathcal{C}} \Delta(c, c')$  and the *relative distance* of the code is  $\min_{c, c' \in \mathcal{C}} \delta(c, c')$ .

There are entire books and courses devoted to error-correcting codes (and more generally, the topic of “coding theory”), so we do not attempt to do the entire field justice here. One of the main reasons codes are useful objects of study is the following observation, which follows by the triangle inequality.

**Observation 3.1.3.** If  $\mathcal{C}$  has distance  $D$ , then the open Hamming balls of radii  $D/2$  about all  $c \in \mathcal{C}$  are disjoint.

In other words, if an adversary takes a codeword  $c \in \mathcal{C}$  then changes fewer than  $D/2$  of its entries to produce a new  $\tilde{c} \in [q]^\ell$ , then another party can later uniquely “decode”  $\tilde{c}$  to obtain  $c$ . That is, there is a unique  $c \in \mathcal{C}$  that could have been modified in this way to obtain  $\tilde{c}$ . This observation is why such  $\mathcal{C}$  are called error-correcting codes: one can imagine that a user wishes to transmit an  $m$ -bit message  $M$  for  $m = \lfloor \log_2 |\mathcal{C}| \rfloor$  to some other user(s). If all users agree ahead of time on some injection  $\text{Enc} : \{0, 1\}^m \rightarrow \mathcal{C}$ , then the sender can transmit  $\text{Enc}(M)$ . The channel of transmission may then have noise that corrupts what was sent (e.g. scratches on a CD, or static on a phone line), but as long as the channel corrupts fewer than  $D/2$  entries of the encoded message, the recipients can uniquely decode to recover  $M$  precisely.

One other parameter typically of interest when discussing codes is the *rate* of the code. Essentially the rate measures the following: if the original messages and the codewords are over the same alphabet  $[q]$ , then by what multiplicative factor do messages blow up after encoding. More precisely, rate is the inverse of this quantity:  $\mathcal{C}$  has enough codewords to support encoding length- $\log_q |\mathcal{C}|$  messages with alphabet  $[q]$ . Meanwhile, these messages are being blown up to length  $t$ . Thus the rate is defined as  $(\log_q |\mathcal{C}|)/\ell$ . Rates are numbers between 0 and 1 and are a form of measuring the



efficiency of the code from a redundancy perspective: the closer the rate is to 1, the more efficient it is.

For our purposes in proving sketching lower bounds in this chapter, we do not need to explicitly describe any code, but rather just show that good codes exist. We will do so via the *probabilistic method*.<sup>2</sup>

**Theorem 3.1.4.** *For any integers  $q, n > 1$ , there exists a code  $\mathcal{C}$  with  $|\mathcal{C}| = n$  and block length  $\ell = O(q \log n)$  with relative distance  $1 - 6/q$ .*

*Proof.* We pick  $c_i \in [q]^\ell$  uniformly at random, and do so independently for  $i = 1, 2, \dots, n$ . We then wish to show that  $\mathcal{C} = \{c_i\}_{i=1}^n$  has the desired property with positive probability. Look at some particular pair  $c \neq c' \in \mathcal{C}$ . Let  $Y_k$  for  $k = 1, 2, \dots, t$  be an indicator random variable for the event  $c_k = c'_k$ . Then the  $Y_k$  are independent Bernoulli with parameter  $1/q$ . Then  $Y = \sum_{k=1}^t Y_k$  equals  $\ell - \Delta(c, c')$  and has  $\mathbb{E}Y = \ell/q$ . By the upper tail of the Chernoff bound in the regime  $\lambda > 2e - 1$  (see Eq. (1.4) and Remark 1.1.5),

$$\mathbb{P}(Y > 6\ell/q) < \exp(-\Omega(\ell/q)),$$

which is less than  $1/n^2$  by choice of  $t$ . Thus by a union bound over all choices of  $c \neq c' \in \mathcal{C}$ , the probability that there *exists* such a pair with  $\ell - \Delta(c, c') > 6\ell/q$  is strictly less than 1. Thus, the desired code exists.  $\square$

**Corollary 3.1.5.** *For any integer  $n > 0$  and any integers  $\ell, q > 1$  such that  $n = q\ell$ , there exists a subset  $\mathcal{B}_{q,\ell}$  of  $\{0, 1\}^n$  satisfying the following properties:*

1. Every  $c \in \mathcal{B}_{q,\ell}$  has support size  $\ell$ , i.e.  $|\{i : c_i \neq 0\}| = \ell$ .
2. For  $c \neq c' \in \mathcal{B}_{q,\ell}$ ,  $|\{i : c_i = c'_i\}| \leq 6\ell/q$ .
3.  $|\mathcal{B}_{q,\ell}| = \exp(\Omega(\ell/q))$ .

*Proof.* The corollary follows from Theorem 3.1.4 by writing codewords in unary. That is, for each  $c \in \mathcal{C} \subset [q]^\ell$ , we convert it into an element of  $\{0, 1\}^{q\ell}$  by replacing each entry  $c_i$  in  $c$  with an element of  $\{0, 1\}^q$  (by putting a 1 in the  $c_i$ th position and 0's in the other  $q - 1$  positions).  $\square$

Now we are ready to present a lower bound against deterministic but approximate algorithms for the distinct elements problem.

**Theorem 3.1.6.** *Suppose  $\mathcal{A}$  is a deterministic streaming algorithm that always outputs a value  $\tilde{t}$  when queried such that  $t \leq \tilde{t} \leq 1.9t$ , where  $t$  is the number of distinct elements. Then  $\mathcal{A}$  uses at least  $cn$  bits of memory for some constant  $c > 0$ .*

*Proof.* We show that the existence of  $\mathcal{A}$  implies the existence of an injection  $\text{Enc} : \mathcal{B}_{q,\ell} \rightarrow \{0, 1\}^S$  for  $q = 100$  and  $\ell = n/q$ , where  $\mathcal{B}_{q,\ell}$  is as in Corollary 3.1.5. Therefore  $S \geq \log |\mathcal{B}_{q,\ell}| = \Omega(n)$ . The injection is defined as follows. Given some  $x \in \mathcal{B}_{q,\ell}$ , we artificially create a stream that contains in some fixed (say sorted) order all  $i$  such that  $x_i = 1$ . We then run the algorithm  $\mathcal{A}$  on it, then define  $\text{Enc}(x)$  to be the memory contents of  $\mathcal{A}$ .

To show that  $\text{Enc}(x)$  is an injection, we show that we can invert it (or ‘‘Decode’’) as in the proof of Theorem 3.1.1. If  $M$  is an  $S$  bit string, then  $\text{Dec}(M)$  is defined as follows:

<sup>2</sup>Briefly, this method works as follows: if one wants to show that an object with some property  $P$  exists, then pick an object randomly from some distribution and show that the picked object has property  $P$  with nonzero probability.

```

for  $c \in \mathcal{B}_{q,\ell}$ :
   $\mathcal{A}.\text{init}(M)$  // initialize  $\mathcal{A}$ 's memory to  $M$ 
  for  $i = 1, 2, \dots, n$ :
    if  $c_i = 1$ :
       $\mathcal{A}.\text{update}(i)$ 
  if  $\mathcal{A}.\text{query}() \leq 1.9\ell$ 
    return  $x$ 

```

Each time through the for loop,  $\mathcal{A}$ 's memory is reset to the contents immediately after the encoder processed  $x$ . Thus the exact number of distinct elements is  $\ell$ , so the reported value to a query will be at most  $1.9\ell$  if  $c = x$ . Meanwhile if  $c \neq x$ , thinking of  $c, x$  as sets (containing the elements corresponding to 1 bits), item (2) of [Corollary 3.1.5](#) implies the true number of distinct elements is  $|c \cup x| = |c| + |x| - |c \cap x| > 2\ell - 6\ell/q = 1.94\ell$ . Thus the output of query will not be at most  $1.9\ell$ . □

The 1.1-approximation in [Theorem 3.1.6](#) can be changed to any approximation factor strictly less than 2 by adjusting  $q$ . It is even possible to show that  $\alpha$ -approximation to the number of distinct elements for any constant  $\alpha > 1$  requires  $\Omega(n/\alpha)$  bits of memory deterministically [[CK16](#)], though we will not present that argument here.

We next show that exact, randomized algorithms must also use  $\Omega(n)$  bits of memory.

**Theorem 3.1.7.** *Suppose  $\mathcal{A}$  is a randomized streaming algorithm that outputs the exact number of distinct elements with success probability at least  $2/3$  for the last query in any fixed sequence of stream updates and queries. Then  $\mathcal{A}$  uses at least  $cn$  bits of memory for some constant  $c > 0$ .*

*Proof.* We first remark that the existence of  $\mathcal{A}$  using space  $S$  and succeeding with probability  $2/3$  implies the existence of  $\mathcal{A}'$  that outputs the exact right answer with failure probability at most  $10^{-6}$  and uses space  $O(S)$ :  $\mathcal{A}'$  simply runs  $O(1)$  copies of  $\mathcal{A}$  in parallel with independent randomness and returns the median query result across all parallel runs as its own query response. The claimed failure probability holds via the Chernoff-Hoeffding bound. We thus assume that  $\mathcal{A}$  in fact fails with probability at most  $10^{-6}$ .

Consider again the set  $\mathcal{B}_{q,\ell}$  from [Corollary 3.1.5](#) for  $q = 100$ . We will show the existence of such  $\mathcal{A}$  implies the existence of an injection  $\text{Enc} : \mathcal{B} \rightarrow \{0, 1\}^S$  for some  $\mathcal{B} \subseteq \mathcal{B}_{q,\ell}$  with  $|\mathcal{B}| \geq |\mathcal{B}_{q,\ell}|/2$ . Thus  $S \geq \log |\mathcal{B}| = \log |\mathcal{B}_{q,\ell}| - 1 = \Omega(n)$ .

We define  $\text{Enc}'$  identically as  $\text{Enc}$  in the proof of [Theorem 3.1.1](#), and  $\text{Dec}'$  is similar, though instead of returning  $x$  we return the closest element (in Hamming) distance from  $\mathcal{B}_{q,\ell}$  to  $x$ . This is not yet the desired injection and inverse because they are randomized procedures and not actual functions, since  $\mathcal{A}$  is a randomized algorithm. For fixed  $x \in \mathcal{B}_{q,\ell}$ , let  $x'$  be the (random) vector recovered by  $\text{Dec}'$ . Let  $Y_{x,i}$  be an indicator random variable for the event that the  $i$ th bit was recovered incorrectly, i.e.  $x'_i \neq x_i$ . Let  $Y_x = \sum_{i=1}^n Y_{x,i}$ . Then  $\mathbb{E}Y_x \leq n/10^6$  by linearity of expectation. Let  $Z_x$  be an indicator random variable for the event  $Y_x > 2n/10^6$ . Then  $\mathbb{P}(Z_x = 1) < 1/2$  by Markov's inequality. Thus by linearity of expectation,

$$\mathbb{E} \sum_{x \in \mathcal{B}_{q,\ell}} Z_x \geq \frac{1}{2} \cdot |\mathcal{B}_{q,\ell}|. \quad (3.1)$$

This expectation is over the randomness of  $\mathcal{A}$ . One can view any randomized algorithm  $\mathcal{A}(x)$  as simply a deterministic algorithm  $\mathcal{A}(x, r)$ , where  $r$  is the (random) string that sources  $\mathcal{A}$  with

randomness. Eq. (3.1) implies that there *exists* a fixed string  $r^*$  such that  $\text{Dec}'(\text{Enc}'(x))$  returns  $x$  for at least half the  $x \in \mathcal{B}_{q,\ell}$  when using  $r^*$  as the source of randomness for  $\mathcal{A}$ . We denote this set of  $x$  where the scheme succeeds as  $\mathcal{B}$ . We then define  $\text{Enc}, \text{Dec}$  as  $\text{Enc}', \text{Dec}'$  but using  $r^*$  as the random string.  $\square$

### 3.1.2 Quantiles

Recall in the quantiles problem, we would like to answer  $\varepsilon$ -approximate rank/select queries after seeing a stream of comparable items. As mentioned, Cormode and Yi showed that any deterministic, comparison-based algorithm for this problem requires  $\Omega(\varepsilon^{-1} \log(\varepsilon n))$  words of memory, showing that the GK sketch is optimal. We will not show their proof here, but rather a much simpler lower bound.

**Theorem 3.1.8.** *Suppose  $\varepsilon \in (10/n, 1)$ . Then any deterministic, comparison-based algorithm for  $\varepsilon$ -approximate rank/select requires at least  $\Omega(1/\varepsilon)$  words of memory.*

*Proof.* Suppose the algorithm stores  $S$  elements of the stream. Including the smallest and largest elements of the stream would increase the space to at most  $S + 2$ . If  $S + 2 \geq n/2$ , then are done. Otherwise, call these elements  $x_1 < x_2 < \dots < x_{S+2}$ . Now any other element in the stream other than those stored is in the interval  $(x_i, x_{i+1})$  for some unique  $i$ ; there are at most  $S - 1$  intervals of this form. Thus these at most  $S - 1$  intervals contain at least  $n - S - 2$  items of the stream, so at least one interval must contain at least  $(n - S - 2)/(S - 1) \geq 2n/S$  stream items. Every item in this interval has identical comparisons with all the stored elements and thus all have the same response to a rank query. Meanwhile, the smallest and largest elements in this interval have a difference of at least  $2n/S - 1$  in rank, which can only be at most  $2\varepsilon n$  if  $S = \Omega(1/\varepsilon)$ . Noting that for two numbers (their ranks) differing by more than  $\ell$  there is integer within  $\ell/2$  of both finishes the proof.  $\square$

We next show a stronger lower bound using a compression-based argument similar to that for distinct elements, based on error-correcting codes. The following lower bounds shows an  $\Omega(1/\varepsilon)$  lower bound even for randomized algorithms, and even when the elements come from a bounded universe  $\{1, \dots, U\}$ .

**Theorem 3.1.9.** *In what follows,  $C, c > 0$  are some universal constants. Let  $\varepsilon \in (0, 1)$  be given, and suppose  $n, U > 1$  are given integers with  $1/\varepsilon > C \cdot \max\{1/n, 1/U\}$ . . Suppose  $\mathcal{A}$  is a randomized streaming algorithm for  $c\varepsilon$ -approximate rank/select with success probability at least  $2/3$  for the last query in any fixed sequence of stream updates and queries, over streams of length  $n$  where elements are in the universe  $\{1, \dots, U\}$ . Then  $\mathcal{A}$  uses at least  $\Omega(\varepsilon^{-1} \log(\varepsilon U))$  bits of memory.*

*Proof.* As in the proof of Theorem 3.1.6, we can assume wlog that  $\mathcal{A}$  actually succeeds with probability at least  $7/8$  while only increasing its space usage  $S$  by at most a constant factor.

Consider all sequences  $S$  of  $1/\varepsilon$  elements of  $[q]$  for  $q = c\varepsilon U$  for some  $c \in (0, 1)$ . Given some sequence  $S$ , we treat it as a length- $(1/\varepsilon)$  string over the alphabet  $[q]$  and encode it with a constant-rate error-correcting code with encoding function  $\text{Enc}$  that can recover from a  $1/4$ -fraction of errors (so the relative distance is more than  $1/2$ ); the length of  $\text{Enc}(S)$  is  $1/\varepsilon' = c^{-1}/\varepsilon$ . The existence of such a code is guaranteed by Theorem 3.1.4 if  $c$  is a small enough constant. Now take the  $i$ th (0-indexed) symbol, call it  $x$ , of  $\text{Enc}(S)$  and change it to  $i \cdot q + x$  for all  $i$ , to get a new string  $S'$  over alphabet  $[U]$  with  $U = q/\varepsilon'$ . Now run  $\mathcal{A}$  on  $S'$  with error parameter  $\varepsilon'/3$  after replicating each letter of  $S'$   $\varepsilon'n$  times (so the stream length is  $n$ ) to get a quantiles data structure. Given the length of  $S'$  and our choice of error parameter, over  $\mathcal{A}$ 's randomness we can recover each symbol

of  $S'$  (and thus  $\text{Enc}(S)$ ) with probability  $7/8$  since  $S'$  is a sorted string and a quantile query for the middle of a repeated character segment of length  $\varepsilon'n$  has a unique correct answer (namely, the repeated character). Thus, the probability we recover fewer than  $3/4$  of the letters is at most  $1/2$ . Thus, there exists a fixed random string to source  $\mathcal{A}$ 's randomness which allows  $\mathcal{A}$  to correctly recover a  $3/4$ -fraction of the letters of  $\text{Enc}(S)$  (and thus decode to recover  $S$ ) for at least half the possible sequences  $S$ . In other words, if the space (in bits)  $\mathcal{A}$  uses is  $s$ , we have an injection from a collection of size at least  $(1/2)(c\varepsilon U)^{1/\varepsilon}$  to bitsrings of length  $s$ , implying  $s = \Omega(\varepsilon^{-1} \log(\varepsilon U))$ .  $\square$

## 3.2 Communication Complexity

Another common method to prove lower bounds is via reductions from problems in communication complexity (though the methods to prove that the communication problems being reduced from are themselves often compression-based). One can imagine representing solving a streaming problem with  $\ell$  updates as a communication game, in which there are  $\ell$  players  $A_1, \dots, A_\ell$  where player  $A_i$  holds update  $i$ . One can imagine that if the players had a streaming algorithm, player  $A_1$  could run  $\mathcal{A}$  on the first update, send  $\mathcal{A}$ 's memory contents as a message to player  $A_2$ , who can then run  $\mathcal{A}$  initialized with that memory on the second update, etc., until player  $A_\ell$  can finally run  $\mathcal{A}$  on their update before calling  $\mathcal{A}.\text{query}()$ . In this way, the memory complexity of solving some streaming problem is captured by the maximum message length required to solve this corresponding communication problem. This translation is not perfect however, since in the communication model we allow players to send arbitrary functions of their input and received message whereas a streaming algorithm must be an algorithm, i.e. it cannot compute uncomputable functions to determine how to update its memory state. For example, nothing prohibits a player in a communication game from deciding what message to send the next player based on the solution to some instance of the halting problem, or some other undecidable problem; its decision on what to send (or player  $A_\ell$ 's decision on what to output) is allowed to be the result of an arbitrarily complex function. As we will see in this section, it is often the case that we can prove optimal memory lower bounds for streaming problems by considering only 2 players, call them Alice and Bob. For many problems we can imagine that Alice's input corresponds to say the first half of stream updates, and Bob holds the second half, and their task is to compute some function on the concatenated stream. Even though Alice and Bob have quite a lot of power (they each know half the entire stream!), it turns out that for many streaming problems solving this associated 2-player communication problem is just as hard.

**Various models of communication games.** Before we continue, we describe a few different types of communication games studied. We describe all these models in the 2-player setting where Alice has some  $x \in \mathcal{X}$ , Bob has some input  $y \in \mathcal{Y}$ , and they would like to send messages back and forth to compute  $f(x, y)$ . Specifically Alice sends a message, to which Bob responds, to which Alice may respond, etc., until one player declares they know the answer and outputs  $f(x, y)$ . We note that this back-and-forth communication corresponds to multiple passes over the input for a streaming algorithm: for example if Alice talks, then Bob talks, then Alice talks again, this corresponds to the algorithm making two passes over the input stream. For a communication protocol  $\Pi$  and inputs  $x, y$ , we define  $\Pi(x, y)$  to be the transcript of all messages sent, and  $|\Pi(x, y)|$  as the total number of bits of all messages in this transcript.

- **Deterministic complexity.** As the name suggests, Alice and Bob act completely deterministically. The complexity of a protocol is then  $\max_{x,y} |\Pi(x, y)|$ . We then define  $D(f)$  as the minimum over all correct communication protocols  $\Pi$  of  $\max_{x,y} |\Pi(x, y)|$ .

- **Randomized complexity.** This is again a worst-case notion of complexity (where  $x, y$  are worst-case inputs), but where Alice and Bob act randomly and must only fail with probability at most  $\delta$ . We define  $R_\delta^{pub}(f)$  as the minimum over all protocols that are correct with probability at least  $2/3$  on any input of  $\max_{x,y} |\Pi(x, y)|$ . The *pub* here denotes that the source of randomness is an infinitely long public random string in the sky that Alice and Bob both have read access to.  $R_\delta^{priv}(f)$  is similarly defined but where Alice and Bob each have their own private source of randomness, which the other player does not know (unless it is communicated, which incurs cost). Clearly  $R_\delta^{priv}(f) \geq R_\delta^{pub}(f)$ , since in the public coin model we can simulate a private coin protocol by having Alice pretend the even-indexed public bits are her private bits, and Bob pretending the odd-indexed public bits are his.
- **Distributional complexity.** In this model  $\mu$  is some distribution over  $(\mathcal{X}, \mathcal{Y})$ , and we assume the inputs are not worst case but rather  $(x, y) \sim \mu$ . Then  $D_\delta^\mu(f)$  denotes the minimum complexity of a protocol that succeeds with probability at least  $1 - \delta$  on  $(x, y)$  drawn from  $\mu$ . Without loss of generality we can assume the protocol is deterministic (since for any randomized protocol, its probability of correctness is the average probability of correctness over its source of randomness, and therefore there must be a fixed random string to source it with that performs at least as well as this average).

We also sometimes consider the “one-way model” in which Alice only sends a single message to Bob, and Bob must then output the answer (Bob never speaks to Alice). For one way complexities, we use the notation  $D^\rightarrow(f), R_\delta^{priv, \rightarrow}(f), R_\delta^{pub, \rightarrow}(f), D_\delta^{\mu, \rightarrow}(f)$ .

We now state a few theorems known in the communication complexity literature without proof. The interested reader can find proofs in the textbook of Kushilevitz and Nisan [KN97].

**Theorem 3.2.1** (Yao’s minimax principle). *For any  $f$  and any  $\delta \in (0, 1)$ ,  $R_\delta^{pub}(f) = \sup_\mu D_\delta^\mu(f)$ .*

**Theorem 3.2.2.** *For any  $f$  and fixed constant  $\delta \in (0, 1/2)$ ,  $R_\delta^{priv} = \Omega(\log(D(f)))$ .*

**Theorem 3.2.3** (Newman’s theorem). *For any  $f : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}$  and  $\epsilon, \delta \in (0, 1)$ ,  $R_{\epsilon+\delta}^{priv} \leq R_\delta^{pub}(f) + O(\log n + \log(1/\epsilon))$ .*

### 3.2.1 Equality

The equality function  $\text{EQUALITY}_n : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}$  is defined by  $\text{EQUALITY}_n(x, y) = 1$  iff  $x = y$ . The following known theorem states that amongst deterministic communication protocols, the trivial protocol of Alice sending Bob her entire input is best possible.

**Theorem 3.2.4.** *For all  $n \geq 1$ ,  $D(\text{EQUALITY}_n) \geq n$ .*

The above theorem implies [Theorem 3.1.1](#) as a corollary.

**Corollary 3.2.5.** *Suppose  $\mathcal{A}$  is a deterministic streaming algorithm that computes the number of distinct elements exactly. Then  $\mathcal{A}$  uses at least  $n$  bits of memory.*

*Proof.* Given such an  $\mathcal{A}$  using  $S$  bits of memory, we define a communication protocol using  $S$  bits of communication for  $\text{EQUALITY}$ . Define  $X = \{i : x_i = 1\}$  and similarly for  $Y$ . Alice runs  $\mathcal{A}$  on the stream consisting of all  $i \in X$  then sends  $\mathcal{A}$ ’s memory contents to Bob as a message. Bob then runs  $\mathcal{A}.\text{query}()$  and immediately outputs 0 if the result is not equal to  $|Y|$ . If he does not reject, he continues running  $\mathcal{A}$  from the memory state sent on all  $i$  in  $Y$  then runs  $\mathcal{A}.\text{query}()$  again. He then outputs 1 iff this second query equals  $|Y|$ . Note the output of this second query is  $|X \cup Y|$ , which is strictly larger than  $|Y|$  iff  $Y \neq X$  since  $X, Y$  are the same size.  $\square$

It is interesting to note that  $\text{EQUALITY}_n$  is extremal for both [Theorems 3.2.2](#) and [3.2.3](#).

**Theorem 3.2.6.**  $R_{1/3}^{\text{priv}}(\text{EQUALITY}_n) = O(\log n)$  and  $R_{1/3}^{\text{pub}}(\text{EQUALITY}_n) = O(1)$

*Proof.* For the private coin model, Alice and Bob each treat their inputs as binary representations of integers in the range  $\{0, \dots, 2^n - 1\}$ . Alice then uses private randomness to pick a uniformly random prime  $p$  in the interval  $[1, 10n^3]$  and sends to Bob the two numbers  $p$  and  $x \bmod p$ . Bob then outputs 1 iff  $x \equiv y \pmod p$ . The total communication is  $O(\log n)$  bits since  $p$  and  $x \bmod p$  are each at most  $n^3$  and thus requires  $3 \log_2 n$  bits each to transmit. As for correctness, this protocol clearly is correct if  $x = y$ . If they are different, then Bob errs iff  $p$  divides  $|x - y|$ . But  $|x - y| < 2^n$  and thus has at most  $n$  prime divisors (note if  $N$  has prime divisors  $p_1, \dots, p_k$  then  $N \geq \prod_{j=1}^k p_j \geq 2^k$ , implying  $k \leq \log_2 N$ ). Meanwhile, the interval  $[1, 10n^3]$  has at least  $(1 - o(1))10n^3 / \ln(10n^3)$  primes by the prime number theorem, and thus the probability that the chosen  $p$  is one of these divisors is at most  $n / ((1 - o(1))10n^3 / \ln(10n^3)) < 1/3$ .

For the public coin model, Alice and Bob use public randomness to agree on two independent uniformly random strings  $r, r' \in \{0, 1\}^n$ . Alice then sends Bob  $\langle r, x \rangle \bmod 2$  and  $\langle r', x \rangle \bmod 2$  (two bits). Bob then outputs 1 iff  $\langle r, x \rangle = \langle r, y \rangle \bmod 2$  and  $\langle r', x \rangle = \langle r', y \rangle \bmod 2$ . Again, this protocol succeeds with probability 1 if  $x = y$ . If they are not equal, then they succeed iff  $\langle x - y, r \rangle = 0 \pmod 2$ . But since they are not equal,  $x - y$  has at least one non-zero entry  $(x - y)_i$ . Then  $\langle x - y, r \rangle = (x_i - y_i)r_i + \sum_{j \neq i} (x_j - y_j)r_j$ . For any conditioning of this sum,  $(x - y)r_i$  either flips the sum or keeps it the same modulo 2 with equal probability. Thus  $\mathbb{P}(\langle x - y, r \rangle = 0) = 1/2$ . Thus Bob fails with probability exactly  $(1/2)^2 = 1/4 < 1/3$ .  $\square$

### 3.2.2 Disjointness

The disjointness problem is defined by  $\text{DISJ}_n(S, T) = 1$  if  $S \cap T = \emptyset$  and 0 otherwise, where  $S, T \subseteq [n]$ . The following theorem is known.

**Theorem 3.2.7.** [[KS92](#), [Raz92](#)]  $R_{1/3}^{\text{pub}}(\text{DISJ}_n) = \Omega(n)$ .

Now consider the  $\ell_\infty$ -norm approximation problem in the streaming model: we see a sequence of integers  $i_1, \dots, i_\ell \in \{1, \dots, n\}$  in a stream, which define a histogram  $x \in \mathbb{R}^n$  where  $x_i$  represents the number of times  $i$  appeared in the stream. Upon a query, we should output some  $\tilde{z}$  such that  $\|x\|_\infty \leq \tilde{z} \leq (1 + \varepsilon)\|x\|_\infty$  with probability at least  $2/3$ . Recall  $\|x\|_\infty = \max_i |x_i|$ .

**Theorem 3.2.8.** Suppose  $\mathcal{A}$  is a randomized streaming algorithm for the  $\ell_\infty$ -norm approximation problem which outputs  $\tilde{z}$  such that  $\|x\|_\infty \leq \tilde{z} \leq 1.1\|x\|_\infty$  with probability at least  $2/3$ . Then  $\mathcal{A}$  uses  $\Omega(n)$  bits of memory.

*Proof.*  $\square$

One might ask: what about approximation factors  $\alpha \gg 1$ ? That is, we should output  $\tilde{z}$  such that  $\|x\|_\infty \leq \tilde{z} \leq \alpha\|x\|_\infty$ . For this case, the work [[BJKS04](#)] considered a  $t$ -player generalization  $\text{DISJ}_{n,t}$ . There are players  $A_1, \dots, A_t$  where  $A_i$  receives as input some  $S_i \subseteq [n]$ . They are promised that either (1)  $A_i \cap A_j = \emptyset$  for all  $i \neq j$ , or (2)  $A_i \cap A_j = \{x\}$  for some  $x \in [n]$  for all  $i \neq j$ . They must decide which case they are in. That work considered the one-way complexity of this problem in which  $A_1$  sends a message to  $A_2$ , who then sends a message to  $A_3$ , etc., until  $A_t$  finally has to output the answer. The complexity measures are for the total amount of bits sent by all players combined.

**Theorem 3.2.9.** [BJKS04]  $R_{1/3}^{\text{pub}}(\text{DISJ}_{n,t}) = \Omega(n/t)$ . In particular, some player must send a message of length at least  $\Omega(n/t^2)$  bits.

**Corollary 3.2.10.** Suppose  $\mathcal{A}$  is a randomized streaming algorithm that provides a better than  $\alpha$ -approximation for the  $\ell_\infty$ -norm approximation problem for some  $\alpha \geq 2$ , with probability at least  $2/3$ . Then  $\mathcal{A}$  uses  $\Omega(n/\alpha^2)$  bits of memory.

*Proof.* The proof is identical to [Theorem 3.2.8](#) but where we reduce from  $\text{DISJ}_{n, \lceil \alpha \rceil}$ . In the case that all sets are pairwise disjoint, we have  $\|x\|_\infty \leq 1$  and so the algorithm will report an answer less than  $\alpha$ . Meanwhile in case (2),  $\|x\|_\infty \geq t = \lceil \alpha \rceil$  and thus will output an answer that is at least  $\lceil \alpha \rceil$ . Thus the two cases are distinguishable by the streaming algorithm, implying some player sends a message of length at least  $\Omega(n/\alpha^2)$  bits. Since each player's message length is exactly the memory usage of  $\mathcal{A}$ , the claim follows.  $\square$

One can also apply the above reasoning to obtain a lower bound for approximating the “ $p$ th moment”  $\|x\|_p^p$  for any  $p \geq 2$ , where  $\|x\|_p^p := \sum_i |x_i|^p$ .

**Corollary 3.2.11.** Suppose  $\mathcal{A}$  is a randomized streaming algorithm that provides a 2-approximation for the  $p$ th moment approximation problem, with probability at least  $2/3$ . Then  $\mathcal{A}$  uses  $\Omega(n^{1-2/p})$  bits of memory.

*Proof.* We reduce from  $\text{DISJ}_{n,t}$  for  $t = \lceil (3n)^{1/p} \rceil$ . In the case the sets are disjoint, we have  $\|x\|_p^p \leq n$ . Meanwhile if some element is in all the pairwise intersections then the  $p$ th moment is at least  $t^p = 3n$ . Thus a 2-approximation can distinguish these cases, solving disjointness.  $\square$

It is known that this  $n^{1-2/p}$  is optimal up to  $\log n$  factors [IW05].

### 3.2.3 Indexing, GapHamming, and Distinct Elements

We wrap up this chapter by showing a chain of reductions that imply  $(1 + \varepsilon)$ -approximation of the number of distinct elements in a data stream requires  $\Omega(1/\varepsilon^2)$  bits of memory (as long as  $\varepsilon > 1/\sqrt{n}$ ). We first introduce the communication problems  $\text{INDEX}_n$  and  $\text{GAPHAMMING}_n$  then show that the former reduces to the latter, which then reduces to distinct elements in the streaming model.

In the  $\text{INDEX}_n$  problem Alice gets  $x \in \{0, 1\}^n$  and Bob gets  $j \in [n]$ , and Bob must output  $x_j$ . We consider one-way model, in which Bob must output the answer based on a single message to Bob. Intuitively this problem is hard, since Alice does not know  $j$  and thus seemingly must tell Bob her entire string for him to succeed with good probability. Our goal is to show  $R_{1/3}^{\text{pub}, \rightarrow}(\text{INDEX}_n) = \Omega(n)$ . We will use [Theorem 3.2.1](#), which states that it suffices to lower bound  $D_{1/3}^{\mu, \rightarrow}(\text{INDEX}_n)$  for some  $\mu$  that we may choose freely. We consider  $\mu$  being the uniform distribution over  $(x, j)$ . Since these are now random variables, we henceforth refer to them as  $(X, J)$ .

Before continuing, we state a few standard results from information theory.

**Lemma 3.2.12** (Chain rule).  $H(Y|X) = H(X, Y) - H(X)$ .

**Lemma 3.2.13** (Fano's inequality). Suppose  $X$  is a random variable finitely supported in  $\mathcal{X}$ . Let  $\hat{X} \stackrel{\text{def}}{=} g(Y)$  be the predicted value of  $X$  for  $g$  a deterministic function also taking values in  $\mathcal{X}$ . Then if  $\mathbb{P}(\hat{X} \neq X) = \delta$ ,

$$H(X|Y) \leq H(X|\hat{X}) \leq H_2(\delta) + \delta \log_2(|\mathcal{X}| - 1),$$

where  $H_2(\delta) \stackrel{\text{def}}{=} -\delta \log_2 \delta - (1 - \delta) \log_2(1 - \delta)$ .

*Proof.* Let  $E$  be an indicator r.v. for  $X \neq \hat{X}$ . Then by the chain rule,

$$H(E, X|\hat{X}) = H(X|\hat{X}) + H(E|X, \hat{X})$$

but also

$$H(E, X|\hat{X}) = H(E|\hat{X}) + H(X|E, \hat{X}) \leq H(E) + H(X|E, \hat{X}) = H_2(\delta) + H(X|E, \hat{X}).$$

Also  $H(E|X, \hat{X}) = 0$ . Thus  $H(X|\hat{X}) \leq H_2(\delta) + H(X|E, \hat{X})$ . But also

$$\begin{aligned} H(X|E, \hat{X}) &= (1 - \delta) \cdot H(X|E = 0, \hat{X}) + \delta \cdot H(X|E = 1, \hat{X}) \\ &= \delta \cdot H(X|E = 1, \hat{X}) \\ &\leq \delta \cdot \log_2(|\mathcal{X}| - 1). \end{aligned}$$

The last inequality is because if  $E = 1$  then  $X \neq \hat{X}$ , so  $X$  can take on one of at most  $|\mathcal{X}| - 1$  values.  $\square$

**Theorem 3.2.14.** For any  $n \geq 1$ ,  $D_\delta^{\text{uniform}, \rightarrow}(\text{INDEX}_n) \geq (1 - H_2(\delta))n$ .

*Proof.* Let  $\Pi(X)$  be the message Alice sends to Bob on input  $X$ . Then given  $\Pi(X)$ ,  $J$ , Bob outputs  $\hat{X}_J$  s.t.  $\mathbb{P}(\hat{X}_J \neq X_J) \leq \delta$ . Thus by Fano's inequality,  $H(X_J|\Pi(X), J) \leq H_2(\delta)$ . By definition of conditional entropy and the chain rule,

$$\begin{aligned} H(X_J|\Pi(X), J) &= \sum_{j=1}^n \mathbb{P}(J = j) \cdot H(X_J|\Pi(X), J = j) \\ &= \frac{1}{n} \sum_{j=1}^n H(X_j|\Pi(X)) \\ &\geq \frac{1}{n} \sum_{j=1}^n H(X_j|\Pi(X), X_1, \dots, X_{j-1}) \\ &= \frac{1}{n} \sum_{j=1}^n H(X_1, \dots, X_j, \Pi(X)) - H(X_1, \dots, X_{j-1}, \Pi(X)) \\ &= 1 - \frac{1}{n} H(\Pi(X)) \\ &\geq 1 - \frac{1}{n} |\Pi|. \end{aligned}$$

The theorem follows by rearranging the inequality  $1 - |\Pi|/n \leq H_2(\delta)$ .  $\square$

In  $\text{GAPHAM}_n$ , Alice and Bob receive  $x, y \in \{0, 1\}^n$ , respectively. Recall the Hamming distance  $\Delta(x, y) := |\{i : x_i \neq y_i\}|$ . In Gap Hamming, Alice and Bob are promised that  $\Delta(x, y)$  is *either* at least  $n/2 + \sqrt{n}$  or at most  $n/2 - \sqrt{n}$  and must decide which.

The following theorem is originally due to [Woo04], but the simpler proof we present here via reduction from  $\text{INDEX}$  was given later in [TSJ08]. The lower bound being restricted to one-way protocols was removed later in [CR12] (see also later simpler proofs [She12, Vid12]).

**Theorem 3.2.15.**  $R_{1/3}^{\text{pub}, \rightarrow}(\text{GAPHAM}_N) = \Omega(N)$ .



*Proof.* We reduce from  $\text{INDEX}_n$  to  $\text{GAPHAM}_N$  for some odd integer  $n = \beta N$  ( $\beta$  a constant to be determined later). Recall in the indexing problem Alice receives  $x \in \{0, 1\}^n$  and Bob receives  $i \in [n]$  and would like to compute  $x_i$  after one single message from Alice. First, Alice forms a vector  $x' \in \{-1, 1\}^n$  where  $x'_i = -1$  if  $x_i = 1$ , and  $x'_i = 1$  if  $x_i = 0$ . The two players will use public randomness to agree upon  $N$  independent, uniformly random vectors  $r^1, \dots, r^N \in \{-1, 1\}^n$ . Alice will then create a vector  $a \in \{-1, 1\}^n$  where  $a_k = \text{sign}(\langle x', r^k \rangle)$ , and Bob creates  $b \in \{-1, 1\}^n$  with  $b_k = r_i^k$ . Note the argument to  $\text{sign}$  can never be 0 since  $n$  is odd. Alice and Bob then let their answer be the answer to  $\text{GAPHAM}_N$  for their inputs  $a, b$ .

We now must argue correctness. We write  $a_k = r_i^k x'_i + (\sum_{j \neq i} r_j^k x'_j) = r_i^k x'_i + \alpha$ . Note if  $\alpha \neq 0$  then  $|\alpha| \geq 2$  since it has an even number of summands, each of which is odd. Thus when  $\alpha \neq 0$ ,  $\text{sign}(a_k) = \text{sign}(\alpha)$ , which is uniformly random in  $\{-1, 1\}$  since  $\alpha$  is a symmetric random variable. Meanwhile if  $\alpha = 0$ , then  $r_i^k x'_i$  equals  $r_i^k$  iff  $x_i = 0$ . Thus, letting  $[[\mathcal{E}]]$  denote 1 if  $\mathcal{E}$  is true and 0 otherwise,

$$\begin{aligned} \mathbb{P}(a_k \neq b_k) &= \mathbb{P}(\alpha \neq 0) \cdot \mathbb{P}(a_k \neq b_k | \alpha \neq 0) + \mathbb{P}(\alpha = 0) \cdot \mathbb{P}(a_k \neq b_k | \alpha = 0) \\ &= \left(1 - \frac{c}{\sqrt{n}}\right) \cdot \frac{1}{2} + \frac{c}{\sqrt{n}} [[x_i = 0]], \end{aligned} \quad (\text{Stirling approximation})$$

which is either  $1/2 - c/\sqrt{n}$  or  $1/2 + c/\sqrt{n}$  depending on whether  $x_i = 0$ . Thus  $\mathbb{E} \Delta(a, b)$  is either  $N/2 - cN/\sqrt{n}$  or  $N/2 + cN/\sqrt{n}$ . By setting  $n = N/(100c^2)$ ,  $\mathbb{E} \Delta(a, b)$  is either  $N/2 + 10\sqrt{N}$  or  $N/2 - 10\sqrt{N}$ . Thus by a Chernoff bound, it is either at least  $N/2 + \sqrt{N}$  or at most  $N/2 - \sqrt{N}$  with large constant probability, which is then decided by the Gap Hamming protocol.  $\square$

**Corollary 3.2.16.** *For any  $\varepsilon \in (1/\sqrt{n}, 1)$ , any algorithm  $\mathcal{A}$  providing a  $(1 + \varepsilon)$ -approximation of the number of distinct elements in a stream with probability at least  $2/3$  must use  $\Omega(1/\varepsilon^2)$  bits of memory.*

*Proof.* We reduce from  $\text{GAPHAM}_N$  for  $N = 1/(5\varepsilon^2)$ . Alice runs  $\mathcal{A}$  on the  $i$  such that  $x_i = 1$  and sends the memory contents of  $\mathcal{A}$  to Bob, as well as  $|\text{support}(x)|$ , i.e.  $|\{i : x_i \neq 0\}|$  (the same as  $|x|$  if one treats  $x$  as a set). Bob then continues running the algorithm on  $i$  such that  $y_i = 1$ . Note the players can do this since although the distinct elements universe is  $[n]$ , we have  $N \leq n$  by the definition of  $N$  and restriction  $\varepsilon > 1/\sqrt{n}$ . Treating  $x, y$  as sets, we have that the number of distinct elements is  $|x \cup y| = \Delta(x, y) + |x \cap y| = \Delta(x, y) + |x| + |y| - |x \cup y|$ . Rearranging,  $\Delta(x, y) = 2|x \cup y| - |x| - |y|$ . Bob is told  $|x|$ , and he knows  $|y|$ . He also can run  $\mathcal{A}.\text{query}()$  to obtain  $2(1 \pm \varepsilon)|x \cup y|$ , which is equal to  $2|x \cup y|$  with an additive error of at most  $2\varepsilon N = 1/\varepsilon < \sqrt{N}$ . Thus this additive error is enough to decide Hamming distance at most  $N/2 - \sqrt{N}$  or at least  $N/2 = \sqrt{N}$  for Gap Hamming.  $\square$



# Chapter 4

## Linear Sketching

The focus of this chapter will be *linear sketching*. This is a general technique for sketching a high-dimensional vector  $x \in \mathbb{R}^n$  where we store  $\Pi x$  in memory for  $\Pi \in \mathbb{R}^{m \times n}$  for some  $m \ll n$ . If  $\Pi$  has a succinct representation (so that we spend  $\ll mn$  space to store it but there is a low-memory algorithm to compute  $\Pi_{i,j}$  given  $i, j$ ), then our sketch reduces the representation of  $x$  from  $n$  units of memory to  $m \ll n$ . This technique is not only commonly used in the design of streaming algorithms, but also when designing distributed algorithms as well as we will discuss later in this chapter.

Regarding streaming algorithms, up until this point we have focused on streaming algorithms in the so-called *insertion-only model*. Specifically, consider the scenario of a vector  $x \in \mathbb{R}^n$  with  $n$  large, and  $x$  starting as the 0 vector. Then we see a sequence of updates  $(i, \Delta)$  each causing the change  $x_i \leftarrow x_i + \Delta$ . For example consider  $\Delta = 1$  for every update, so that  $x$  is simply a histogram tracking the number of occurrences of each item from some size- $n$  universe in the stream

In the above setup, there are three popularly studied models:

1. **insertion-only:** Each update has  $\Delta = 1$  (the assumption in previous chapters).
2. **strict turnstile:** Some updates  $\Delta$  can be negative, but we are given the promise that  $\forall i, x_i \geq 0$  at all times. This makes sense for example in graph streaming, in which case  $n = \binom{|V|}{2}$  for vertex set  $V$ , and  $x_e$  represents the multiplicity of edge  $e$ . In a dynamic graph edges may be inserted and deleted, but the multiplicity of an edge would never be negative.
3. **general turnstile:** Anything goes. Updates can be negative, and entries in  $x$  can be negative as well. For example, we may wish to process one time period's updates with  $\Delta = -1$  and another period with  $\Delta = +1$  so that we can later query  $x$  being the difference vector of the two time periods.

In this chapter, when discussing streaming we focus on algorithms for the strict and general turnstile models. Note that if we store  $y = \Pi x$  in memory, the update  $(i, \Delta)$  causes the change  $x \leftarrow x + \Delta \cdot e_i$ , and thus  $y \leftarrow y + \Delta \cdot \Pi^i$ , where  $\Pi^i$  is the  $i$ th column of  $\Pi$  (since generally  $Az = \sum_i z_i A^i$  for any matrix-vector product  $Az$ ). All known algorithms for both strict and general turnstile are actually linear sketches. It is in fact known [LNW14, AHLW16] that *any* algorithm in these two models can be converted into a linear sketch with only a logarithmic factor loss in space complexity, if the algorithm is required to be correct on very long streams<sup>1</sup>.

---

<sup>1</sup>When this assumption does not hold, in certain cases one can do better than linear sketching [KP20].

## 4.1 Heavy hitters

In this section we discuss two types of (related) problems, referred to as *heavy hitters* and *point query*. For both of these problems, we consider  $x$  being updated in the turnstile streaming model. There is also an integer parameter  $k > 0$  given at the beginning of the stream.

In the  $\ell_1$  *point query* problem with parameter  $k$  (we sometimes say  $(k, \ell_1)$ -*point query*), we must answer queries of the form  $\text{query}(i)$ , for  $i \in [n]$ , with a value in the range  $x_i \pm \|x\|_1/k$ . Here  $k$  is a parameter known at the beginning of the stream.

In  $\ell_1$  *heavy hitters* with parameter  $k$  (or  $(k, \ell_1)$ -*heavy hitters*), there is only one query, and we must answer it with a set  $L \subset [n]$  such that

1.  $|L| = O(k)$
2.  $|x_i| > \|x\|_1/k \implies i \in L$

The indices  $i$  satisfying the second criterion are called *k-heavy hitters* or *k-frequent* (sometimes we drop the  $k$  if clear from context).

Note  $\|x\|_1/k = (1/k) \sum_{i=1}^n |x_i|$ . The number of  $i$  which can be strictly larger than an  $\alpha$ -fraction of this sum must be strictly less than  $1/\alpha$ , and thus there can be at most  $k - 1$  heavy hitters. We are thus requiring that the list  $L$  being returned not be more than a constant factor larger than the absolute biggest it need to be to contain all the heavy hitters.

**Remark 4.1.1.** The tail versions of these problems are also commonly studied. Define  $x_{\text{tail}(k)}$  to be the vector  $x$  but with the top  $k$  entries in magnitude all zeroed out. Then in the  $\ell_1$ -tail point query problem, we would like error  $\|x_{\text{tail}(k)}\|_1/k$  instead of  $\|x\|_1/k$ . Similarly for heavy hitters, we require that  $|x_i| > \|x_{\text{tail}(k)}\|_1/k \implies i \in L$ . Note under this definition, there can be at most  $2k - 1$  heavy hitters. Specifically, ignoring the  $k$  “head” indices (those not in the tail), there are fewer than  $k$  tail indices that can contribute strictly more than a  $1/k$  fraction of the total tail mass.

The following lemma shows why point query and heavy hitters are related.

**Lemma 4.1.2.** *Suppose there is an algorithm  $\mathcal{A}$  solving  $(3k, \ell_1)$ -point query with failure probability at most  $\delta/n$  and using  $S$  words of memory. Then there is an algorithm  $\mathcal{A}'$  solving  $(k, \ell_1)$ -heavy hitters with failure probability at most  $\delta$  and using space at most  $S + 3k$ .*

*Proof.* Algorithm  $\mathcal{A}'$  uses  $\mathcal{A}$  to process the stream. Then to answer a heavy hitters query, it loops through all  $i \in [n]$  and point queries each one, remembering the  $4k$  indices  $i$  with the point query values (breaking ties arbitrarily). Henceforth condition on the event that all  $n$  point queries return correct values, which happens with failure probability at most  $\delta$  by the union bound. Note then that any  $k$ -frequent index  $i$  will have a point query value strictly larger than  $\|x\|_1/k - \|x\|_1/(3k) = (2/3)\|x\|_1/k$ . Meanwhile, any index  $i$  with  $|x_i| < \|x\|_1/(3k)$  will have a point query value strictly less than  $(2/3)\|x\|_1/k$ , and thus will not appear larger than any actual  $k$ -frequent item. Since there are at most  $3k$  indices satisfying  $|x_i| \geq \|x\|_1/(3k)$ , our return list is thus guaranteed to contain all  $k$ -frequent indices.  $\square$

### 4.1.1 CountMin sketch

We here describe the CountMin sketch [CM05], which solves  $\ell_1$  point query in the general turnstile model. We will describe it here in the strict turnstile model. We now describe the operation of the algorithm:

1. We store hash functions  $h_1, \dots, h_L : [n] \rightarrow [t]$ , each chosen independently from a 2-wise independent family.
2. We store counters  $C_{a,b}$  for  $a \in [L]$ ,  $b \in [B]$  with  $B = 2k$ ,  $L = \lceil \log_2(1/\delta) \rceil$ .
3. Upon an update  $(i, \Delta)$ , we add  $\Delta$  to all counters  $C_{a, h_a(i)}$  for  $a = 1, \dots, L$ .
4. To answer  $query(i)$ , we output  $\min_{1 \leq a \leq L} C_{a, h_a(i)}$ .

Note that our total memory consumption is  $O(L)$  to store the seeds that specify all  $L$  hash functions, as well as  $O(BL)$  words to store the counters  $C_{a,b}$ . Thus the total memory consumption is  $O(BL)$  words.

**Lemma 4.1.3.** *CountMin.query(i) returns  $x_i \pm \|x\|_1/k$  w.p.  $\geq 1 - \delta$ .*

*Proof.* Fix  $i$ , let  $Z_j = 1$  if  $h_r(j) = h_r(i)$ ,  $Z_j = 0$  otherwise. Now note that for any  $r \in [L]$ ,  $C_{r, h_r(i)} = x_i + \sum_{j \neq i} x_j Z_j := x_i + E$ . We have  $\mathbb{E} E = \sum_{j \neq i} |x_j| \cdot \mathbb{E} Z_j = \sum_{j \neq i} |x_j|/B \leq \|x\|_1/(2k)$ . Thus by Markov's inequality,  $\mathbb{P}(E > \|x\|_1/k) < 1/2$ . Thus by independence of the  $L$  rows of the CountMin sketch,  $\mathbb{P}(\min_r C_{r, h_r(i)} > x_i + \|x\|_1/k) < 1/2^L \leq \delta$ .  $\square$

Thus we easily obtain the following theorem via [Lemma 4.1.2](#).

**Theorem 4.1.4.** *There is an algorithm solving the  $\ell_1$   $k$ -heavy hitter problem in the strict turnstile model with failure probability  $\delta$ , space  $O(k \log(n/\delta))$ , update time  $O(\log(n/\delta))$ , and query time  $O(n \log(n/\delta))$ .*

It is also possible to obtain the tail guarantee for point query from the CountMin sketch with the same memory up to a constant factor.

**Lemma 4.1.5.** *The CountMin sketch as above but with  $B \geq 4k$  guarantees that for any  $i$ ,  $query(i) = x_i \pm \|x_{tail(k)}\|_1/k$  w.p.  $\geq 1 - \delta$ .*

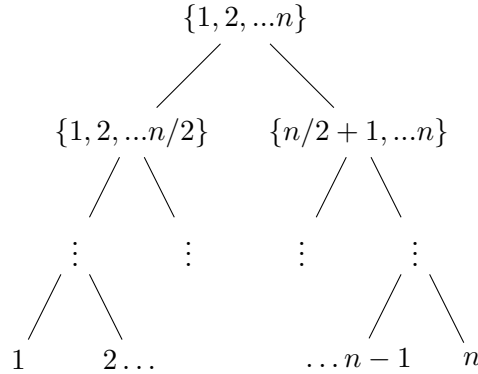
*Proof.* We let  $H \subset [n]$  denote the indices of the top  $k$  entries in magnitude in  $x$ , and  $T := [n] \setminus H$  be the remaining indices (the support of the tail). For fixed  $r \in [L]$  we write  $C_{r, h_r(i)} - x_i = \sum_{j \in H \setminus \{i\}} |x_j| Z_j + \sum_{j \in T \setminus \{i\}} |x_j| Z_j := E_1 + E_2$  where  $Z_j$  is an indicator random variable for the event  $h_r(i) = h_r(j)$ . Consider the event  $\mathcal{F}_r$  that  $h_r(i) \notin h_r(H \setminus \{i\})$ . This happens with probability at least  $1 - |H|/B \geq 3/4$  by a union bound over all  $j \in H \setminus \{i\}$  (recall  $|H| = k$ ). When  $\mathcal{F}_r$  occurs,  $E_1 = 0$ . Now consider the event  $\mathcal{F}'_r$  that  $|E_2| \leq \|x_T\|_1/k = \|x_{tail(k)}\|_1/k$ , where  $x_S$  denotes the projection of  $x$  onto  $S \subset [n]$  (i.e. zero out all entries that are not in  $S$ ). By an analysis similar to the proof of [Lemma 4.1.3](#),  $\mathbb{P}(\mathcal{F}'_r) \geq 3/4$ . Thus by a union bound  $\mathbb{P}(\mathcal{F}_r \wedge \mathcal{F}'_r) \geq 1/2$ , i.e. for each  $r$   $C_{r, h_r(i)}$  provides at most the desired error with probability at least  $1/2$ . Since each  $C_{r, h_r(i)}$  is either equal to  $x_i$  or an overestimate,  $\min_r C_{r, h_r(i)}$  fails to give the desired error guarantee iff each  $C_{r, h_r(i)}$ , which happens with probability at most  $2^{-L} \leq \delta$ .  $\square$

**Remark 4.1.6.** Similar guarantees for point query and heavy hitters are obtained in the general turnstile model, but with a slightly larger  $B$  and  $L$  (by constant factors) and with the estimator returned as the median of the  $C_{r, h_r(i)}$  instead of the minimum. For example, if we pick  $B = 3k$  then we have  $\mathbb{P}(|C_{r, h_r(i)} - x_i| > \|x\|_1/k) < 1/3$ . The Chernoff bound thus implies that the median estimator succeeds with probability  $1 - \delta$  for  $L = C \log(1/\delta)$  for sufficiently large constant  $C > 0$ .

### Speeding up query time

While the above algorithm gives *some* correct heavy hitter algorithm with small space, the query time is quite slow. Here we discuss the *dyadic trick* technique of [CM05] to speed up query.

Consider a perfect binary tree whose leaves are in correspondence with  $[n]$ .



There are  $1 + \lg n$  levels of the tree. We imagine level  $j$  of the tree (with the root being level 0) corresponds to a  $2^j$ -dimensional vector  $x^{(j)}$  being updated. Each node in the tree is a coordinate of the vector at the corresponding level, and the value at that coordinate is the sum of the two values of the children (with  $i$ th leaf simply having value  $x_i$ ). Then when we see an update  $(i, \Delta)$ , we imagine that this update happens to all coordinates that are ancestors of the  $i$ th leaf. Then what we actually store in memory is  $1 + \lg n$  CM sketches, one per level. Then upon an update, we feed that update to the appropriate coordinate at the CM sketch at every level.

Each CM sketch is instantiated to solve the  $4k$ -heavy hitters problem with failure probability  $\eta = \delta/(4k \log n)$ . Thus the total space is  $O(k \log(1/\eta) \log n) = O(k \log((k \log n)/\delta) \log n)$ .

To answer a query, the key insight is that in the strict turnstile model *the value at any ancestor of a node is at least as big as the value at that node*, and furthermore the  $\ell_1$  norm of the implicit vector at each level of the tree is exactly the same. Therefore, if  $i$  is a heavy hitter for the vector  $x$  at the lowest level of the tree, then *every* ancestor of  $i$  is a heavy hitter at its level as well. Since there can only be at most  $2k$  indices that are  $2k$ -heavy hitters, this suggests the following depth first search tree. We move down the tree starting from the root (the root vertex is certainly a heavy hitter for its 1-dimensional vector). At each level  $j$  of the tree, we keep track of a list  $L_j$  of heavy hitters at that level ( $L_j$  should contain all  $k$ -heavy hitters of the vector at its level and not contain any item that is not at least a  $2k$ -heavy hitter). Then, for each of the two children of an index in  $L_j$ , we point query that child using the CM sketch at level  $j + 1$ . If a child has point query output at least  $(3/4)\|x\|_1/k$ , we include it in  $L_{j+1}$  (note that computing  $\|x\|_1$  exactly in the strict turnstile model is trivial: maintain a counter). Finally, our final output list  $L$  is simply the list corresponding to the bottom-most level of the tree.

**Correctness.** Note that since each  $L_j$  has size at most  $2k$ , we point query at most  $4k$  children on the next level. Thus the total number of queries is at most  $4k \log n$  (if it is ever the case that we find ourselves querying more, we can simply output Fail). Thus since we only do at most  $Q \leq (4 \lg n)/\alpha$  queries, since each CM sketch has failure probability at most  $\delta/Q$ , by a union bound the probability that any point query we ever make fails is at most  $\delta$ . Conditioned on such failure not occurring, if some heavy hitter leaf  $i$  is not included in the final  $L$ , look at the highest ancestor  $i'$  of  $i$  on some level  $j$  which was not included in  $L_j$ .  $i'$  cannot be the root, since it is included in  $L_0$ . Thus the

parent of  $i'$  was included in  $L_{j-1}$ , which implies  $i'$  was point queried; this is a contradiction to it not being included, since we conditioned on all point queries succeeding.

**Complexity.** The space used is  $O(k \lg(1/\eta) \log n) = O(k \log((k \log n)/\delta) \log n)$  (in words) as mentioned. The query time is the same. The update time is  $O(\log n \lg(1/\eta)) = O(\log n \lg((k \log n)/\delta))$ .

Though we will not discuss it here, currently the best known algorithm for  $\ell_1$  heavy hitters in the turnstile model is the ExpanderSketch of [LNN16]. It achieves  $O(k \log(n/\delta))$  words of space,  $O(\log(n/\delta))$  update time, and  $O(k \log(n/\delta) \text{poly}(\lg n))$  query time to achieve failure probability  $\delta$ .

### 4.1.2 CountSketch

The CountSketch [CCF04] algorithm is similar to the CountMin sketch, except it provides the guarantee that point query returns an estimate equal to  $x_i \pm \|x\|_2/\sqrt{k}$  (called “ $\ell_2$  point query”). It also solves  $\ell_2$  heavy hitters, which requires returning a list of size  $L$  containing all  $i$  such that  $|x_i| > \|x\|_2/\sqrt{k}$ . One could also study the tail version of this problem, in which one wishes the point query error to be at most  $\|x_{\text{tail}(k)}\|/\sqrt{k}$ , and where  $\ell_2$  tail heavy hitters are defined to be the  $i$  such that  $|x_i| > \|x_{\text{tail}(k)}\|/\sqrt{k}$ .

We first show that obtaining the  $\ell_q$  tail guarantee is always at least as good as obtaining the  $\ell_p$  tail guarantee for  $q > p$ , up to potentially changing the  $k$  by a factor of 2. Thus an  $\ell_2$  tail guarantee is better than an  $\ell_1$  tail guarantee. For example, consider the vector  $x = (\sqrt{n}, 1, 1, \dots, 1)$ . Then index  $i = 1$  is a 1-heavy hitter in  $\ell_2$  even in the non-tail version of the problem, whereas it is only  $O(1/\sqrt{n})$ -heavy in  $\ell_1$ .

**Lemma 4.1.7.** *For any  $1 \leq p < q$ ,  $\|x_{\text{tail}(2k)}\|_q/k^{1/q} \leq \|x_{\text{tail}(k)}\|_p/k^{1/p}$ .*

*Proof.* For simplicity assume  $x$  is infinite-dimensional (pad it with zeroes). Let  $S_1 \subset [n]$  of size  $k$  be the set of  $i$  with the largest  $|x_i|$  values (break ties arbitrarily). Let  $S_2$  be the indices of the next  $k$  largest entries, etc. Then

$$\begin{aligned}
\frac{1}{k^{1/q}} \|x_{\text{tail}(2k)}\|_q &= \frac{1}{k^{1/q}} \left( \sum_{j=3}^{\infty} \|x_{S_j}\|_q^q \right)^{1/q} \\
&\leq \frac{1}{k^{1/q}} \left( \sum_{j=3}^{\infty} k \cdot \|x_{S_j}\|_{\infty}^q \right)^{1/q} \\
&\leq \frac{1}{k^{1/q}} \left( \sum_{j=3}^{\infty} k \cdot \frac{\|x_{S_{j-1}}\|_p^q}{k^{q/p}} \right)^{1/q} && (\forall i \in S_j, |x_i| \leq (\|x_{S_{j-1}}\|_p^p/k)^{1/p}) \\
&= \frac{1}{k^{1/p}} \left( \sum_{j=3}^{\infty} \|x_{S_{j-1}}\|_p^q \right)^{1/q} \\
&= \frac{1}{k^{1/p}} \left( \sum_{j=3}^{\infty} \|x_{S_{j-1}}\|_p^p \right)^{1/p} && (\forall v, \|v\|_q \leq \|v\|_p) \\
&= \frac{1}{k^{1/p}} \|x_{\text{tail}(k)}\|_p
\end{aligned}$$

□

The CountSketch works similarly to the CountMin sketch, but differs slightly in the following way: we in addition select  $L$  hash functions  $h_1, \dots, h_L : [n] \rightarrow \{-1, 1\}$  independently from a 2-wise independent family. Then to process an update  $(i, \Delta)$ , we add  $\sigma_r(i)\Delta$  to  $C_{r, h_r(i)}$  for  $r = 1, \dots, L$  (instead of only adding  $\Delta$  as in the CountMin sketch). For point query, we pick  $B = 9k$  and  $L = C \log(1/\delta)$  for some constant  $C > 0$ . To answer a query, we output the median of  $\sigma_r(i)C_{r, h_r(i)}$  over all  $r \in [L]$ .

**Lemma 4.1.8.** *CountSketch.query(i) returns  $x_i \pm \|x\|_2/\sqrt{k}$  w.p.  $\geq 1 - \delta$ .*

*Proof.* Fix  $r$  and let  $Z_j$  be an indicator random variable for the event  $h_r(i) = h_r(j)$ . Write the error random variable  $E_r := \sigma_r(i)C_{r, h_r(i)} - x_i = \sum_{j \neq i} \sigma_r(i)\sigma_r(j)Z_j x_j$ . Then the lemma is equivalent to showing that the median over  $r$  of  $|E_r|$  is at most  $\|x\|_2/\sqrt{k}$  with probability at least  $1 - \delta$ . For this to hold, by the Chernoff bound and choice of  $L$  it suffices to show that for any  $r \in [L]$ ,  $\mathbb{P}(|E_r| > \|x\|_2/\sqrt{k}) < 1/3$ .

$$\begin{aligned}
\mathbb{E}|E_r| &\leq \sqrt{\mathbb{E}E_r^2} && \text{(Lemma 1.1.9)} \\
&= \left( \mathbb{E} \left[ \sum_{j, j' \neq i} \sigma_r(j)\sigma_r(j')Z_j Z_j' x_j x_{j'} \right] \right)^{1/2} \\
&= \left( \mathbb{E} \left[ \sum_{j \neq i} Z_j x_j^2 + \sum_{\substack{j, j' \neq i \\ j \neq j'}} \sigma_r(j)\sigma_r(j')Z_j Z_j' x_j x_{j'} \right] \right)^{1/2} \\
&= \left( \sum_{j \neq i} (\mathbb{E} Z_j) x_j^2 + \sum_{\substack{j, j' \neq i \\ j \neq j'}} (\mathbb{E} \sigma_r(j)\sigma_r(j')) (\mathbb{E} Z_j Z_j') x_j x_{j'} \right)^{1/2} && \text{(linearity of expectation)} \\
&= \left( \sum_{j \neq i} (\mathbb{E} Z_j) x_j^2 + \sum_{\substack{j, j' \neq i \\ j \neq j'}} (\mathbb{E} \sigma_r(j)) (\mathbb{E} \sigma_r(j')) (\mathbb{E} Z_j Z_j') x_j x_{j'} \right)^{1/2} && \text{(2-wise independence)} \\
&= \left( \sum_{j \neq i} (\mathbb{E} Z_j) x_j^2 \right)^{1/2} \\
&\leq \frac{\|x\|_2}{\sqrt{B}} \\
&= \frac{1}{3} \cdot \frac{\|x\|_2}{\sqrt{k}}
\end{aligned}$$

Thus by Markov's inequality,  $\mathbb{P}(|E_r| > \|x\|_2/\sqrt{k}) < 1/3$ , as desired.  $\square$

The following corollary then holds by combining Lemma 4.1.8 with Lemma 4.1.2.

**Corollary 4.1.9.** *The CountSketch solves  $\ell_2$   $k$ -heavy hitters in general turnstile streams with failure probability  $\delta$  using  $O(k \log(n/\delta))$  words of memory.*

The following lemma also holds by combining Lemma 4.1.8 with an analysis almost identical to that of Lemma 4.1.5.



**Lemma 4.1.10.** *The CountSketch sketch as above but with  $B \geq 15k$  guarantees that for any  $i$ ,  $\text{query}(i)$  returns a value  $x_i \pm \|x_{\text{tail}(k)}\|_2/\sqrt{k}$  w.p.  $\geq 1 - \delta$ .*

## 4.2 Graph sketching

In this section we show the perhaps surprising result that linear sketching can be used to solve combinatorial problems on graphs. Specifically, we consider a dynamic (multi)graph on  $n$  vertices in which vertices can be both inserted and deleted. This model can be faithfully represented in the strict turnstile model, where  $x$  has dimension  $\binom{n}{2}$ ;  $x_e$  specifies the presence (or number of copies) of edge  $e$  in the graph. An insertion of edge  $e$  then corresponds to the turnstile update  $(e, +1)$ , whereas an edge deletion corresponds to the update  $(e, -1)$ . Solving graph problems naively would require remembering  $x$  exactly, i.e.  $\Omega(n^2)$  memory in the worst case. In this section we describe the ‘‘AGM sketch’’ [AGM12], which shows that dynamic spanning forest can be solved using  $O(n \log^c n)$  bits of memory with success probability  $1 - 1/\text{poly}(n)$ . The AGM sketch achieves the exponent with  $c = 3$ , though in these notes we describe a simpler version achieving  $c = 4$ . One can also obtain bounds based on the failure probability as a separate parameter  $\delta$ ; see [NY19]. Note that being able to query the dynamic spanning forest allows for solving many other problems, such as finding an  $s$ - $t$  path, global connectivity, or  $s$ - $t$  connectivity. It is furthermore known that the  $\Omega(n \log^3 n)$  bound is tight; any algorithm that reports a spanning forest with probability at least even 1% must use this much memory [NY19]. Yu recently showed a stronger lower bound in the distributed sketching model, in which all vertices and a ‘‘referee’’ share public randomness each vertex knows only its own neighborhood and must send a short message to the referee; he showed that even if the query is not required to output an entire spanning forest but simply a single bit indicating whether the graph is connected, the average message length must be  $\Omega(\log^3 n)$  bits [Yu21] (the AGM sketch provides a matching upper bound for spanning forest in this model as well). After the paper [AGM12], linear sketching has been proven useful for a wide variety of dynamic graph problems; see [McG14] for a survey.

The AGM sketch uses certain other data structures as subroutines, which we describe first. Both these data structures operate in the general turnstile model.

### 4.2.1 $k$ -sparse recovery

Recall  $\text{support}(x) \subseteq [n]$  denotes the indices  $i$  such that  $x_i \neq 0$ . The parameter  $k$  is given at the beginning of the stream, and there is a single query. The answer to a query is as follows: if  $|\text{support}(x)| \leq k$ , the query simply returns  $x$  exactly (the indices in the support together with their values); otherwise, the query response can be arbitrary. We show that this problem can be solved deterministically using  $O(k \log(nM))$  bits of memory if we are promised that all update amounts are integers and no entry of  $x$  is ever larger than  $M$  in magnitude. Note with this restriction it suffices to solve the problem over  $\mathbb{F}_p$  for any prime  $p > M$ , as any  $x_i$  will equal  $x_i \bmod p$ . The solution we discuss will also require  $p \geq n$ , and thus we will work over any prime  $p > \max\{M, n\}$ . The total space will be at most  $2k \lceil \log p \rceil$  bits.

Recall in linear sketching we maintain  $\Pi x$  in memory for some  $\Pi \in \mathbb{R}^{m \times n}$  (in this case  $\Pi \in \mathbb{F}_p^{m \times n}$ ).  $k$ -sparse recovery is possible iff  $\Pi x \neq \Pi y$  for  $x, y$  distinct  $k$ -sparse vectors, which is equivalent to  $\Pi(x - y) \neq 0$ . Noting  $x - y$  is  $2k$ -sparse, this requirement is thus equivalent to  $\Pi z \neq 0$  for any  $2k$ -sparse vector  $z$ . If  $S$  denotes the support of  $z$ , then  $\Pi z = \Pi_S z$ , where  $\Pi_S$  is the  $m \times |S|$  submatrix of  $\Pi$  keeping only the columns indexed by  $S$ . Thus our requirement for  $\Pi$  is that all of its  $m \times 2k$  submatrices have full column rank. We will pick  $m = 2k$ , so these are square submatrices; thus having full column rank is equivalent to  $\det(\Pi_S) \neq 0$  for all  $2k \times 2k$  submatrices

of  $\Pi$ . We will specifically pick  $\Pi$  to be the transpose of a Vandermonde matrix. Specifically, pick  $x_1 \neq x_2 \neq \dots \neq x_n \in \mathbb{F}_p$  (this is why we require  $p \geq n$ , to guarantee at least  $n$  distinct elements in  $\mathbb{F}_p$ ) and set  $\Pi_{i,j} = x_j^{i-1} \pmod p$  for  $i, j \in [n]$ . For concreteness, we could pick  $x_j = j$ . The following known fact, which we will not prove here, implies that any  $2k \times 2k$  submatrix of  $\Pi$  has nonzero determinant.

**Fact 4.2.1.** *Let  $A \in F^{r \times r}$  be such that  $A_{i,j} = x_j^{i-1}$  for  $i, j \in [r]$  for some field  $F$ . Then*

$$\det(A) = \prod_{1 \leq i < j \leq n} (x_i - x_j)$$

The above fact is usually written for  $A^\top$ , but note  $\det(A) = \det(A^\top)$  for any  $A$ . Note [Fact 4.2.1](#) implies  $\det(A) \neq 0$  if the  $x_i$  are distinct.

**Lemma 4.2.2.** *Suppose  $A \in F^{m \times n}$  is such that every  $m \times 2k$  submatrix of  $A$  has full column rank, where  $F$  is a field. Then there is an algorithm running in  $\binom{n}{2k} \cdot \text{poly}(n)$  to recover  $x$  given  $y = Ax$  for any  $k$ -sparse  $x$ .*

*Proof.* We loop over all  $S \subset [n]$  of size exactly  $2k$  (our guess for a set containing the support of  $x$ ) and compute  $x' = \Pi_S^{-1}y$ . If  $x'$  is  $k$ -sparse, then we form  $x$  as the  $n$ -dimensional vector  $x$  with  $x_S = x'$  and  $x_{[n] \setminus S} = \vec{0}$  and return  $x$ .  $\square$

**Remark 4.2.3.** For the particular scheme we propose, it is possible to actually recover  $x$  from  $\Pi x$  in  $O(k^2 \text{polylog}(p))$  time via an algorithm called *syndrome decoding*, though we will not cover it here.

## 4.2.2 SupportFind

In this problem, we would like a randomized data structure which, if  $x = 0$ , reports `null`. Otherwise, if  $x \neq 0$  it should return some  $i \in \text{support}(x)$  with probability at least  $1 - \delta$  (with probability  $\delta$  it is allowed to behave arbitrarily). There are no promises regarding the vector  $x$ ; it may or may not be sparse, yet the query algorithm should still succeed.

We describe an algorithm, the JST sketch, for this problem due to [\[JST11\]](#) which uses  $O(\log(1/\delta) \log^2 n)$  bits of memory if all entries in  $x$  are promised to be integers which are at most  $\text{poly}(n)$  in magnitude. It is known that this bound is optimal even if the entries of  $x$  are promised to always be either 0 or 1 [\[KNP<sup>+</sup>17\]](#).

The JST sketch uses the geometric sampling technique of [Subsection 2.2.3](#). Specifically, we pick a hash function  $h : [n] \rightarrow [\log_2 n]$  with  $\mathbb{P}(h(i) = j) = 1/2^j$  (other than for  $j = \log_2 n$ , in which case we have  $\mathbb{P}(h(i) = j) = 1/2^{j-1}$  so that the probabilities add to one). For now we assume  $h$  is a perfectly random hash function, though we discuss in [Remark 4.2.5](#) how this can be relaxed using bounded independence. For each  $j \in [\log_2 n]$ , we also instantiate a  $k$ -sparse recovery data structure  $A_j$  from [Subsection 4.2.1](#) with  $k = C \log(1/\delta)$  for a sufficiently large constant  $C > 0$  and with  $M = \text{poly}(n)$ .

To process `update( $i, \Delta$ )`, we simply call  $A_{h(i)}.update(i, \Delta)$ . To process a query, we loop from  $j = \log_2 n$  down to 1 and for each such  $j$  call  $A_j.query()$ . For the first (i.e. largest) value of  $j$  for which the query is not the zero vector, we return any index in the support of the query response. See [Subsection 4.2.2](#) for pseudocode.

**Theorem 4.2.4.** *If  $x = 0$ , `null` is returned with probability 1. Otherwise, the probability some  $i \in \text{support}(x)$  is returned is at least  $1 - \delta$ .*

```

for  $j = \log_2 n, \dots, 1$ :
   $z \leftarrow A_j.\text{query}()$ 
  if  $z \neq 0$ :
    return any  $i$  such that  $z_i \neq 0$ 
return null //  $x$  is the zero vector

```

*Proof.* The case  $x = 0$  is clear, as all  $A_j$  will return the zero vector when queried. Also clear is the case  $|\text{support}(x)| \leq k := C \log(1/\delta)$ , since every  $A_j$  will receive a vector with support size at most that of  $x$  (and thus will return its received vector exactly), and at least one of the  $A_j$  must receive a nonzero vector if  $x \neq 0$  since every index of  $x$  is hashed to exactly one  $A_j$ .

For the remainder of the proof we thus focus on the case that  $\text{support}(x) \geq k$ . Let  $t$  denote  $|\text{support}(x)|$ . Let  $x^{(j)}$  denote the vector  $x$  where we zero out all coordinates  $i$  such that  $h(i) \neq j$ , and define the random variable  $T_j := |\text{support}(x^{(j)})|$ . For some (large) constant  $c$  such that  $1 < c < C$ , let  $j^* \in [\log_2 n]$  be such that  $c \log(1/\delta) \leq t/2^{j^*} < 2c \log(1/\delta)$ . Such  $j^*$  must exist since  $t > C \log(1/\delta)$ . We define two events:  $\mathcal{E}_1$  is the event  $\max_{j \geq j^*} T_j \leq k$ , and  $\mathcal{E}_2$  is the event  $T_{j^*} \geq 1$ . Note if both events occur, then our query output is guaranteed to be correct. This is because  $\mathcal{E}_1$  implies  $A_j.\text{query}()$  will correctly return  $x^{(j)}$  for all  $j \geq j^*$ , and  $\mathcal{E}_2$  implies that at least one of these  $x^{(j)}$  is nonzero (since in particular  $x^{(j^*)} \neq 0$ ). We then show  $\mathbb{P}(\neg \mathcal{E}_1 \vee \neg \mathcal{E}_2) \leq \delta$ , by the union bound.

We bound  $\mathbb{P}(\neg \mathcal{E}_1)$  itself by a union bound over  $j \geq j^*$ . Note  $\mathbb{E} T_j = t/2^j$ . Then  $\mathbb{P}(T_j > k) = \mathbb{P}(T_j > (k2^j/t) \cdot \mathbb{E} T_j) < (k2^j/t)^{-C'k}$  (see Eq. (1.7)). Summing over  $j \geq j^*$  gives a geometric series dominated by its largest term, which is the term for  $j^*$ , which is  $(k2^{j^*}/t)^{-C'k} \leq (k/(c \log(1/\delta)))^{-C'k} = (C/c)^{-C'k} < \delta/2$  for  $C$  sufficiently larger than  $c$ .  $\mathbb{P}(\neg \mathcal{E}_2)$  is also bounded by the Chernoff bound. We have  $\mathbb{E} T_{j^*} \in [c \log(1/\delta), 2c \log(1/\delta)]$ . Thus  $\mathbb{P}(\neg \mathcal{E}_2) = \mathbb{P}(T_{j^*} = 0)$ , which is at most  $\delta/2$  by the Chernoff bound.  $\square$

**Remark 4.2.5.** It is a useful fact to know that tail bounds imply moment bounds and vice versa. In one direction, if we have a bound on all moments  $\|Z\|_p$  then we have a tail bound via Markov's inequality:  $\mathbb{P}(|Z| > \lambda) < \inf_p \{\lambda^{-p} \|Z\|_p^p\}$  (recall  $\|Z\|_p := (\mathbb{E}|Z|^p)^{1/p}$ ). The value  $p \geq 1$  can be chosen to minimize the right hand side. In the other direction,  $\|Z\|_p^p = \int_0^\infty px^{p-1} \mathbb{P}(|Z| > x) dx$  via integration by parts. Thus a tail bound on  $|Z|$  yields moment bounds. Now, since the Chernoff bound gives strong tail bounds, one can use this correspondence to obtain the implied moment bounds on  $|\sum_i X_i - \mu|$  for all  $p \geq 1$ , and from those moment bounds re-derive the Chernoff bound itself by choosing  $p$  optimally based on  $\lambda$  and  $\mu$ , which is determined by  $p$ -wise independence of the  $X_i$  if  $p$  is an even integer. The punchline is that if one were to carry out this calculation exercise, one would find that whenever the Chernoff bound yields tail probability  $\delta$ , it sufficed to choose  $p = O(\log(1/\delta))$ , so that the  $X_i$  could be  $O(\log(1/\delta))$ -wise independent (see also [BR94, SSS95], which take different approaches to showing this). This observation allows one to select the  $h$  in the JST sketch from an  $O(\log(1/\delta))$ -wise independent family, so that it only takes  $O(\log(1/\delta) \log n)$  bits to represent.

### 4.2.3 AGM sketch

Before describing the AGM sketch, we first design a non-streaming algorithm. Imagine that we proceed in  $R = \log_2 n$  rounds. We start each round with a partition that is a refinement of the partition of vertices into connected components, and in the first round each vertex is in its own partition. Now, at the beginning of each round we ask each partition (which we henceforth call a *super-vertex*) to identify an edge leaving it and entering another partition. At the end of the

vertex we then merge along the set of all identified edges: if two super-vertices are connected by an identified edge, we merge them into an even bigger super-vertex. The spanning forest we return is the set of all identified edges across the rounds. Note this algorithm is correct since the number of non-maximal components at least halves in each round, and we started with  $n$  components and at the end there are at least 1.

We now describe how to implement the above approach in the streaming model. Imagine each vertex  $u$  keeps track of a *signed neighborhood vector*  $x_u \in \mathbb{R}^{\binom{n}{2}}$ . For any edge  $e$  such that either  $e$  is not actually in the graph, or  $e$  does not contain vertex  $u$ , we set  $(x_u)_e = 0$ . However if  $e = (u, v)$  is actually in the graph, we set  $(x_u)_e = 1$  if  $u < v$  or  $(x_u)_e = -1$  if  $u > v$ . The key reason for signing in this way is that if  $A$  is a partition (or super-vertex), then if we define  $x_A := \sum_{u \in A} x_u$ , then  $\text{support}(x_A)$  is exactly the set of edges with exactly one endpoint in  $A$  and one endpoint in a different super-vertex. Recall that the data structure of [Subsection 4.2.2](#) is a (randomized) linear sketch. We pick  $R = \log n$  such linear sketching matrices independently  $\Pi_1, \dots, \Pi_R$  each with failure probability  $\delta < 1/(n^{1+\beta}R)$  (if our desired failure probability is  $1/n^\beta$ ). We store in memory  $\Pi_r x_u$  for all vertices  $u \in [n]$  and all  $r \in [R]$ . The total space is thus  $O(nR \log^3 n) = O(n \log^4 n)$  bits. We simulate the offline algorithm by, in each round  $r$ , forming sketches for each super-vertex  $A$  as  $\sum_{u \in A} \Pi_r x_u$ , which by linearity is just  $\Pi_r(\sum_{u \in A} x_u) = \Pi_r x_A$ . We then query for each  $A$  to obtain edges leaving each super-vertex. The probability that we always obtain correct edges is the probability that no `SupportFind` query ever fails, which is at most  $\delta n R < 1/n^b$  by the union bound.

We remark that one may be tempted to pick only a single  $\Pi$  and store all  $\Pi x_u$  in memory. Then we can reuse the same  $\Pi$  in each round. Doing so unfortunately is incorrect. This is because our guarantees for randomized algorithms  $\mathcal{A}$  are of the following form: for all inputs  $x$ ,  $\mathbb{P}(\mathcal{A} \text{ gives the correct answer on input } x) \geq 1 - \delta$ . Here the probability is over some random string  $\alpha$  that provides *mathcal{A}* with its source of randomness. But by the order of quantifiers, this means  $x$  must be fixed *before* drawing  $\alpha$  randomly. In other words,  $x$  is not allowed to depend on  $\alpha$ . However if we use the same  $\Pi$  in each round, then the fact that we merged certain vertices after round 1 is because of the (random) set of edges our `SupportFind` data structure happened to identify in round 1. Then we form super-vertices  $A$  based on these identified edges, so our next query in round 2, i.e. the fact that we are asking about certain  $A$ , is correlated with the randomness of  $\Pi$ . The AGM sketch avoids this by using fresh random  $\Pi_r$  in each round, independent of the linear sketches used in previous rounds. Thus the queries being asked of  $\Pi_r$  are uncorrelated with the randomness used to specify  $\Pi_r$ .

**Remark 4.2.6.** The version of the AGM sketch described here uses  $O(n \log^4 n)$  bits of memory, though it is possible to implement it using  $O(n \log^3 n)$  bits of memory. The improvement comes from slightly changing the definition of the `SupportFind` problem which the AGM sketch relies on. Instead, imagine defining the problem so that there are two types of failure modes, with separate failure probabilities  $\delta_1$  and  $\delta_2$ . In the first failure mode, the data structure should output `Fail`. In the second failure mode, it may fail without warning (i.e. it simply outputs an index not actually in the support of  $x$ ). Note some failures of the first form are tolerable for the AGM sketch; it simply means there are some supervertices in some rounds which fail to identify an outgoing edge. This is acceptable though, as long as most supervertices in most rounds do identify an outgoing edge (and we can increase  $R$  by a constant factor to compensate). Thus we can set  $\delta_1$  to be some small constant, e.g.  $1/10$ . Failures of the second type though are deadly, since even one such failure causes the entire algorithm to fail. We thus set  $\delta_2 = 1/\text{poly}(n)$ . It is possible to show that such a data structure, with these two failure modes where one is allowed to happen fairly often, can be implemented more memory-efficiently than the structure in [Subsection 4.2.2](#) (though the design of the data structures are very similar), leading to the optimal implementation of the AGM sketch.

For details, see the appendix of [NY19].

### 4.3 Norm estimation

The problem of  $\ell_p$  norm estimation, i.e. providing a multiplicative approximation to  $\|x\|_p := (\sum_i |x_i|^p)^{1/p}$ , was first investigated by Alon, Matias, and Szegedy [AMS99]. Specifically, let  $F_p$  denote the  $p$ th moment  $\|x\|_p^p = \sum_{i=1}^n |x_i|^p$ . As usual, we would like a  $(1 + \varepsilon)$ -approximation to  $F_p$  with probability  $2/3$ . It turns out there is a phase transition in the space complexity of  $F_p$  estimation.

- $0 \leq p \leq 2$ : it is known that  $\text{poly}(\varepsilon^{-1} \log n)$  words of space is achievable [AMS99, Ind06]. For  $p = 0$ , we treat  $0^0$  as 0 and any nonzero element to the 0th power as 1, which is the limit of  $F_p$  for  $p \downarrow 0$ . Note that in insertion-only streams,  $F_0$  is simply the distinct elements problem from Section 2.2.
- For  $p > 2$  and constant  $\varepsilon$ , it is known the space complexity is  $n^{1-2/p}$  up to logarithmic factors [BJKS04, IW05]. That is, there are both upper and lower bounds. Recall this fact was discussed in and below Corollary 3.2.11.

#### 4.3.1 AMS sketch

We now look at the case  $p = 2$ , which is solved by the AMS sketch [AMS99]. Let  $\sigma \in \{-1, 1\}^{m \times n}$  be drawn from a 4-wise independent family for some  $m < n$  to be determined later. We need  $O(\log(mn)) = O(\log n)$  bits, i.e. one machine word, to represent  $\sigma$ . Define  $\Pi \in \mathbb{R}^{m \times n}$  by  $\Pi_{i,j} = \sigma_{i,j}/\sqrt{m}$  so that  $y = \Pi x$  satisfies  $y_i = \sum_{j=1}^n \sigma_{i,j} x_j / \sqrt{m}$ . We estimate  $\|x\|_2^2$  by  $\|\Pi x\|_2^2 = \|y\|_2^2$ .

**Analysis.**

$$\begin{aligned}
 \mathbb{E} y_r^2 &= \frac{1}{m} \mathbb{E} \left( \sum_{j=1}^n \sigma_{r,j} x_j \right)^2 \\
 &= \frac{1}{m} \left[ \|x\|_2^2 + \mathbb{E} \sum_{j \neq j'} \sigma_{r,j} \sigma_{r,j'} x_j x_{j'} \right] \\
 &= \frac{1}{m} \left[ \|x\|_2^2 + \sum_{j \neq j'} (\mathbb{E} \sigma_{r,j} \sigma_{r,j'}) x_j x_{j'} \right] \\
 &= \frac{1}{m} \left[ \|x\|_2^2 + \sum_{j \neq j'} (\mathbb{E} \sigma_{r,j}) (\mathbb{E} \sigma_{r,j'}) x_j x_{j'} \right] \quad (2\text{-wise independence}) \\
 &= \frac{1}{m} \|x\|_2^2,
 \end{aligned}$$

and thus  $\|y\|_2^2 = \sum_{r=1}^m y_r^2$  has expectation  $\|x\|_2^2$ . Next we need to estimate the variance in order to apply Chebyshev's inequality. Observe that

$$\mathbb{E} (\|y\|_2^2 - \mathbb{E} \|y\|_2^2)^2 = \frac{1}{m^2} \mathbb{E} \left( \sum_{r=1}^m \sum_{j \neq j'} \sigma_{r,j} \sigma_{r,j'} x_j x_{j'} \right)^2$$

$$= \frac{1}{m^2} \sum_{r_1, r_2} \sum_{\substack{j_1 \neq j_2 \\ j_3 \neq j_4}} (\mathbb{E} \sigma_{r_1, j_1} \sigma_{r_1, j_2} \sigma_{r_2, j_3} \sigma_{r_2, j_4}) x_{r, j_1} x_{r, j_2} x_{r, j_3} x_{r, j_4} \quad (4.1)$$

$$= \frac{2}{m} \sum_{j_1 \neq j_2} x_{j_1}^2 x_{j_2}^2 \quad (4.2)$$

$$\leq \frac{2}{m} \|x\|_2^4,$$

To see Eq. (4.2), observe that if  $r_1 \neq r_2$  then the four  $\sigma_{r, j}$  in Eq. (4.1) all have different indices and thus by 4-wise independence the expectation is zero. Thus we need only consider the case  $r_1 = r_2 = r$ . In this case, we must either have  $j_1 = j_3, j_2 = j_4$  or  $j_1 = j_4, j_2 = j_3$  else at least one random sign will appear with exponent one and make the expectation zero. Now for a fixed  $j < j'$ , the term  $x_j^2 x_{j'}^2$  appears twice in the summation Eq. (4.2) (once for  $j_1 = j, j_2 = j'$  and once for  $j_1 = j', j_2 = j$ ), whereas it appears four times in Eq. (4.1). We thus multiply by the factor two to compensate.

Thus by Chebyshev's inequality, for  $m \geq 6/\varepsilon^2$  the probability our estimator is outside of  $[(1 - \varepsilon)\|x\|_2^2, (1 + \varepsilon)\|x\|_2^2]$ , i.e. deviates from its expectation by more than  $\varepsilon\|x\|_2^2$ , is at most

$$\frac{2\|x\|_2^4}{m} \cdot \frac{1}{\varepsilon^2\|x\|_2^4} \leq 1/3.$$

**Remark 4.3.1.** One can also obtain the same result by letting  $\Pi$  be the CountSketch matrix (as shown in [TZ12]). That is, we pick random  $h : [n] \rightarrow [m]$  from a 2-wise independent family and  $\sigma \in \{-1, 1\}^n$  from a 4-wise independent family and define  $\Pi \in \mathbb{R}^{m \times n}$  to be the matrix with exactly one nonzero per column:  $\Pi_{h(j), j} = \sigma_j$  for each  $j \in [n]$ . Then one can show  $\mathbb{E} \|\Pi z\|_2^2 = \|z\|_2^2$  and  $\text{Var}[\|\Pi z\|_2^2] = O(1/m)\|z\|_2^4$ , so that by Chebyshev's inequality for  $m = O(1/\varepsilon^2)$  we have  $\|\Pi z\|_2^2$  is a  $(1 \pm \varepsilon)$ -approximation of  $\|z\|_2^2$  with probability at least  $2/3$ .

### 4.3.2 Indyk's $p$ -stable sketch

The AMS sketch gives a memory-efficient sketch for  $p = 2$ , but what about other norms? In [Ind06], Indyk showed that a memory-efficient streaming algorithm exists for estimating  $\ell_p$  norms for any  $p \in (0, 2)$  (when  $p \leq 1$  is not a norm, but  $\|x\|_p := (\sum_i |x_i|^p)^{1/p}$  is still a well-defined function). To accomplish this, he made use of  *$p$ -stable distributions*.

**Definition 4.3.2.** A probability distribution  $\mathcal{D}_p$  over  $\mathbb{R}$  is said to be  *$p$ -stable* if for  $Z, Z_1, \dots, Z_n$  independently drawn from  $\mathcal{D}_p$  and for any fixed  $x \in \mathbb{R}^n$ , the random variable  $\sum_{i=1}^n x_i Z_i$  is equal in distribution to  $\|x\|_p \cdot Z$ .

Some examples are the standard normal distribution  $\mathcal{N}(0, 1)$ , which is 2-stable. This holds since  $x_i g_i$  is a gaussian with variance  $x_i^2$ , and the sum of independent gaussians is a gaussian whose variance is the sum of the individual variances. Another less-known example is the Cauchy distribution, which is 1-stable; it has probability density function  $\varphi(x) = 1/(\pi(1 + x^2))$ . It is a known theorem that such distributions exist iff  $p \in (0, 2]$ . Note that  $p$ -stable random variables for  $p \neq 2$  cannot have bounded variance, since otherwise the sum of independent copies would have to be gaussian as a limiting distribution by the central theorem. In fact, it is known that any  $p$ -stable distribution must have tail bounds  $\mathbb{P}(|Z| > \lambda) = O(1/(1 + \lambda)^p)$  for all  $\lambda > 0$  [No10] (see also [Nel11, Theorem 42]); this implies that such distributions cannot exist for  $p > 2$  (since otherwise they would have bounded variance, violating the central limit theorem). For  $p < 2$  in fact the tail is precisely  $\Theta(1/(1 + \lambda)^p)$ , so they do not have bounded absolute  $q$ th moments for any  $q \geq p$ .

Though  $p$ -stable distributions do not necessarily have closed form density functions, they do have closed form characteristic functions, i.e.  $\hat{\varphi}_Z(t) = \mathbb{E} e^{itZ}$  (i.e. the Fourier transform of the pdf). Namely,  $\hat{\varphi}_Z(t) = e^{|t|^p}$ . Note then for the random variable  $x_i Z_i$ ,  $\hat{\varphi}_{x_i Z_i}(t) = \mathbb{E} e^{i(tx_i)Z_i} = e^{|x_i|^p |t|^p}$ . Since adding two independent random variables convolves their pdfs, it pointwise multiplies their characteristic functions, so that  $\sum_i x_i Z_i$  has characteristic function  $e^{\|x\|_p^p |t|^p}$ .

We first describe an idealized version of Indyk's  $p$ -stable sketch. Pick a matrix  $\Pi \in \mathbb{R}^{m \times n}$  where the  $\Pi_{i,j} = Z_{i,j}$  are i.i.d.  $p$ -stable random variables, scaled so that  $\mathbb{P}(Z \in [-1, 1]) = 1/2$  (note that scaling a  $p$ -stable distribution by a fixed constant preserves  $p$ -stability, i.e. if  $Z$  follows a  $p$ -stable distribution then  $\alpha Z$  follows a  $p$ -stable distribution as well). We maintain  $y = \Pi x$  in memory, and we estimate  $\|x\|_p$  as  $\text{median}_{1 \leq r \leq m} |y_r|$ .

**Analysis.** Let  $I_S : \mathbb{R} \rightarrow \mathbb{R}$  be the indicator of the set  $S$ , i.e.  $I_S(x) = 1$  if  $x \in S$ , and it equals zero otherwise. Then since  $y_r/\|x\|_p$  is  $p$ -stable with scale factor 1, we have  $\mathbb{E} I_{[-1,1]}(y_r/\|x\|_p) = 1/2$  for each  $r \in [m]$ . Because the  $p$ -stable distribution has a pdf which is both bounded and continuous (this is known; see full discussion in [KNW10a]), we also have the following two facts by linearity of expectation:

- $\mathbb{E} \left[ \sum_{r=1}^m I_{[-1-\varepsilon, 1+\varepsilon]} \left( \frac{y_r}{\|x\|_p} \right) \right] = \frac{m}{2} + \Theta(\varepsilon m) \geq \frac{m}{2} + c_1 \varepsilon m.$
- $\mathbb{E} \left[ \sum_{r=1}^m I_{[-1+\varepsilon, 1-\varepsilon]} \left( \frac{y_r}{\|x\|_p} \right) \right] = \frac{m}{2} - \Theta(\varepsilon m) \leq \frac{m}{2} + c_2 \varepsilon m.$

Note that if  $\sum_r I_{[-1-\varepsilon, 1+\varepsilon]}(y_r/\|x\|_p) > m/2$  then strictly more than half the  $y_r$  satisfy  $|y_r| \leq (1+\varepsilon)\|x\|_p$ , and similarly  $\sum_r I_{[-1-\varepsilon, 1+\varepsilon]}(y_r/\|x\|_p)$  implies that strictly less than half the  $y_r$  satisfy  $|y_r| \leq (1-\varepsilon)\|x\|_p$ . Thus if both these events happen simultaneously, we indeed have that the median estimate is  $(1 \pm \varepsilon)\|x\|_p$ . To show that this happens with good probability, we use Chebyshev's inequality. Specifically, for any  $r$  we have  $\text{Var}[I_S(y_r/\|x\|_p)] \leq 1$  since the range of  $I_S$  is  $[0, 1]$ . Thus the variances of the above sums are each at most  $m$  by independence of the  $y_r$ . Chebyshev's inequality and a union bound thus implies that the probability that either sum deviates from its expectation by more than  $3\sqrt{m}$  is at most  $2/9$ . We can then ensure  $3\sqrt{m} < \max\{c_1 \varepsilon m, c_2 \varepsilon m\}$  by picking  $m \geq 9(\min\{c_1, c_2\})^{-2}/\varepsilon^2$ .

Of course, the two main issue with this idealized algorithm, as in [Subsection 2.2.1](#), are precision and independence. The  $\Pi_{i,j}$  are real numbers, but our computer can only perform finite-precision arithmetic; this can be dealt with by simply rounding the  $\Pi_{i,j}$  to appropriate precision before doing computation. We will not delve into those details here as they are fairly routine, and we instead refer the reader to [KNW10a]. Regarding the independence, we used independence in two places: (1) to argue that the variance of the sum of indicators equals the sum of the variances, and (2) to know that  $y_r/\|x\|_p$  is  $p$ -stable, so that we can estimate  $\mathbb{E} I_S(y_r/\|x\|_p)$ . For (1), note this holds even if the random variables summed are only 2-wise independent. Thus we can simply have that the random seeds  $s_1, \dots, s_m$  used to generate the rows of  $\Pi$  are not fully independent, but rather is a simply from a 2-wise independent sample space. For (2), it is known (though unfortunately quite complicated to prove so we will not do so here), that  $k$ -wise independence suffices for  $k = O(1/\varepsilon^p)$ .

**Theorem 4.3.3** ([KNW10a]). *Let  $Z_1, \dots, Z_n$  be i.i.d. from  $\mathcal{D}_p$ , and let  $Y_1, \dots, Y_n$  be  $k$ -wise independent from  $\mathcal{D}_p$ . Then for any fixed  $x \in \mathbb{R}^n$ ,*

$$\sup_{t \in \mathbb{R}} \left| \mathbb{E} I_{(-\infty, t]} \left( \sum_i x_i Z_i \right) - \mathbb{E} I_{(-\infty, t]} \left( \sum_i x_i Y_i \right) \right| < O(1/k^{1/p}).$$

**Theorem 4.3.3** says that the CDFs of the distributions of  $\sum_i x_i Z_i$  and  $\sum_i y_i Z_i$  are close everywhere. Note  $I_{(a,b]}$  is simply  $I_{(\infty,b]} - I_{(\infty,a]}$ , so **Theorem 4.3.3** implies that the amount of probability mass in any interval is the same in the two distributions up to an additive  $O(1/k^{1/p})$ . We have only talked about  $k$ -wise independence for uniform distributions in this course, but note that any distribution can be generated from a uniform random variable in  $[0, 1]$  via the inverse CDF (which in our case we will discretize to finite precision, i.e. integer multiples of  $\gamma$  for some small  $\gamma$ ). Specifically for  $p$ -stable random variables, it is known how to do this generation efficiently [CMS76].

### 4.3.3 Branching programs and pseudorandom generators

Although Indyk’s  $p$ -stable sketch could be derandomized to obtain an optimal algorithm via  $k$ -wise independence, it is unfortunately quite technical to show that this is true, and it was not even known until about a decade after Indyk published his algorithm. It turns out though that there is a simple generic way to derandomize many streaming algorithms, and that is by modeling their execution as that of a *Read-Once Branching Program* (ROBP) then using a generic Pseudorandom Generator (PRG) against such objects, such as Nisan’s PRG [Nis92].

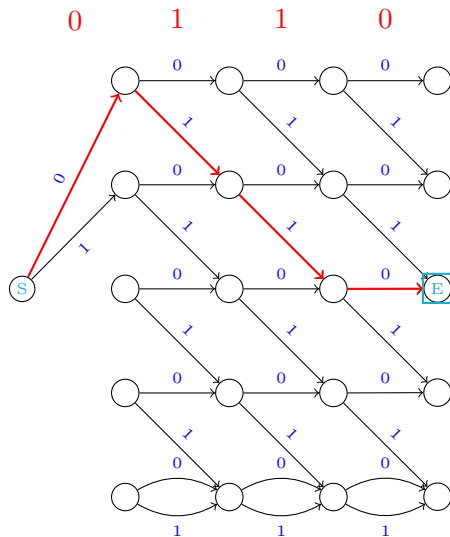


Figure 4.1: ROBP calculating the sum of 4 input bits. The start node is labeled **S**. There are 5 layers (including the start node layer), where the  $k$ th layer (0-indexed) keeps track of the partial sum of the first  $k$  bits. As the partial sum is always between 0 and 4, we thus have 5 nodes per layer. The input bits are “0 1 1 0”, written at the top, and these bits cause us to transition along the edges colored red to land at a particular final vertex (labeled **E**, corresponding to a sum of 2).

**Definition 4.3.4.** A *read-once branching program* is a directed, layered graph with layers  $0, 1, 2, \dots, L$ . Here  $L$  is called the *length*. The 0th layer has just one vertex, called the *start vertex*, and every other layer has the same number of vertices  $W$ , called the *width*. Each vertex, except in layer  $L$ , has out-degree exactly  $d$  for some value of  $d$ , and each edge goes into a vertex in the next layer. The out-edges from a given vertex are also labeled  $0, 1, \dots, d-1$ . Vertices in layer  $L$  have no out-edges. Such an ROBP can be viewed as computing a function  $f : \{0, 1, \dots, d-1\}^L \rightarrow [W]$ , where on input  $(x_1, x_2, \dots, x_L) \in \{0, 1, \dots, d-1\}^L$ , the evaluation of the function is the index of the vertex one



arrives at in the final layer by starting at the start vertex and in the  $i$ th step following the edge out of the current vertex with label  $x_i$ .

For example, Fig. 4.1 describes an ROBP which calculates the sum  $x_1 + x_2 + x_3 + x_4$  of 4 input bits.  $L$  is 4 and the width  $W$  is 5. Essentially each layer is simply keeping track of the running sum of the input bits seen so far (which is a number in  $\{0, 1, 2, \dots, 4\}$  and hence at most 5 states are needed per layer). The left layer with a yellow  $S$  is the start vertex, and the red highlighted path shows the evaluation of  $f$  on input  $(0, 1, 1, 0)$ . The vertex in layer  $L$  boxed in yellow denotes the output state, which here we imagine is indexed as 2 (counting vertices from top to bottom, starting at 0). Note that one can imagine an ROBP of width  $W$  as representing a computation using only  $\lceil \log_2 W \rceil$  bits of memory, since that is the memory required to remember the current state.

In the case of for example Indyk's  $p$ -stable sketch, imagine that we generate each  $\Pi_{i,j}$  from a uniform random number in  $[0, 1]$  with only  $b$  bits of precision, i.e. an integer multiple of  $1/2^b$ , and we round the resulting  $p$ -stable random variable so that it only requires  $b'$  bits of precision to store. Then since the output of the algorithm is invariant under permutation of the stream (since it applies a linear sketch), the algorithm's estimate of  $\|x\|_p$  would be identical if we saw all updates to  $i = 1$  first, then  $i = 2$ , etc. Consider the following computation which sees the entries of  $\Pi$  in row-major order, i.e. we see the entries in the order  $\Pi_{1,1}, \Pi_{1,2}, \dots, \Pi_{1,n}, \Pi_{2,1}, \dots, \Pi_{2,n}, \dots, \Pi_{m,1}, \dots, \Pi_{m,n}$  (more specifically, we see the  $b$ -bit integers that specify each of these entries).

Indyk ROBP:

```

 $c_{low} \leftarrow 0$ 
 $c_{high} \leftarrow 0$ 
for  $r = 1, \dots, m$ :
   $C \leftarrow 0$ 
  for  $j = 1, \dots, n$ :
     $C \leftarrow C + x_j \cdot \Pi_{r,j}$ 
  if  $C \leq (1 - \varepsilon)\|x\|_p$ :
     $c_{low} \leftarrow c_{low} + 1$ 
  if  $C \leq (1 + \varepsilon)\|x\|_p$ :
     $c_{high} \leftarrow c_{high} + 1$ 

```

One can model the evolution of the memory state of Subsection 4.3.3 as a ROBP with  $d = 2^b$ ,  $L = mn$ , and the width sufficiently large to remember  $C, c_{low}, c_{high}$ . If each  $x_j$  is an integer bounded by  $T$  in magnitude, then we can store  $C$  using  $b' + \log(nT)$  bits of precision. Also  $c_{low}, c_{high}$  are always integers in the set  $\{0, 1, \dots, m\}$ . Thus the total space we need to identify our current state (i.e. a vertex in a layer) is  $\lceil \log_2 W \rceil = O(b' + \log(nT) + \log m)$  bits. The final state reveals  $c_{low}, c_{high}$ , which is the number of  $r$  such that  $|y_r| \leq (1 - \varepsilon)\|x\|_p$  and  $|y_r| \leq (1 + \varepsilon)\|x\|_p$ , respectively. As seen in Subsection 4.3.2, with probability  $7/9$  these are simultaneously strictly less than  $m/2$  and strictly more than  $m/2$ .

The goal of a PRG is to preserve the distribution over final states with good probability, so that feeding the ROBP truly random bits to generate the  $\Pi_{i,j}$  versus *pseudorandom* bits results in a distribution over vertices in the final layer that is almost distinguishable.

**Definition 4.3.5.** The *total variation distance* between two probability distributions  $\mathcal{D}, \mathcal{D}'$  is  $\|\mathcal{D} - \mathcal{D}'\|_{TV} := \sup_S |\mathbb{P}_{X \sim \mathcal{D}}(X \in S) - \mathbb{P}_{X' \sim \mathcal{D}'}(X' \in S)|$ .

From our perspective, we can imagine that  $\mathcal{D}$  generates a sequence  $X$  of  $nmb$  independent, uniform bits that specifies the  $\Pi_{i,j}$ , and  $\mathcal{D}'$  generates a pseudorandom sequence  $X'$  of  $nmb$  bits. We can define  $S$  to be the set of inputs  $x$  (i.e. the  $nmb$  bits specifying all of  $\Pi$ ) such that for  $f$  being the ROBP representing Subsection 4.3.3,  $f(x)$  leads to a memory state that implies Indyk's output

is correct, i.e.  $(1 \pm \varepsilon)\|x\|_p$ . Then if  $\mathcal{D}, \mathcal{D}'$  have TV-distance at most  $\epsilon$ , then using the pseudorandom bits generated by  $\mathcal{D}'$  must still lead to a correctness probability of at least  $2/3 - \epsilon$ . Nisan's PRG gives us just that.

**Theorem 4.3.6.** *Let  $\mathcal{U}_{A,t}$  denote the uniform distribution on  $A^t$ . For any  $S, L \geq 1$  there exists a function  $G_{\text{nisan}} : \{0, 1\}^s \rightarrow (\{0, 1\}^S)^L$  for  $s = O(S \log L)$  such that for any  $f$  calculated by an ROBP with width  $2^S$ , degree  $d = 2^S$ , and length  $L$ ,*

$$\|f(\mathcal{U}_{\{0,1\}^s, L}) - f(G_{\text{nisan}}(\mathcal{U}_{\{0,1\}^s}))\|_{TV} \leq \exp(-\Omega(S)).$$

Furthermore, for any  $x \in \{0, 1\}^s$  and any  $j \in [L]$ ,  $(G_{\text{nisan}}(x))_j$  can be computed in space  $O(S \log L)$ .

We can thus derandomize Indyk's algorithm by letting the  $\Pi_{i,j}$  be specified by  $G_{\text{nisan}}$  applied to a short  $s$ -bit random string for  $s = O(S \log L) = O((b' + \log(nT) + \log m) \cdot \log(mn)) = O((b' + \log(nT)) \log n)$  bits since  $m \leq n$ . It can be shown that one can take  $b' = O(\log(nT/\varepsilon))$  [KNW10a], leading to  $O(\log(nT/\varepsilon) \cdot \log n)$  bits overall to specify a sufficiently good pseudorandom matrix  $\Pi$ . Overall this leads to an algorithm for  $\ell_p$  norm estimation using  $O(\varepsilon^{-2} \log(nT/\varepsilon) + \log(nT/\varepsilon) \log n)$  bits of memory (the first summand is for maintaining the actual sketch, and the second is for remembering the seed to Nisan's PRG). Using  $k$ -wise independence instead of Nisan's PRG eliminates the second term, leading to a space-optimal algorithm (matching a lower bound of [KNW10a]), at the cost of requiring a much more complicated analysis.

## Chapter 5

# Johnson-Lindenstrauss Transforms

The following “Johnson-Lindenstrauss lemma” (JL lemma) has been highly impactful in the design of algorithms for high-dimensional data.

**Theorem 5.0.1** (JL lemma [JL84]). *For any  $\varepsilon \in (0, 1)$  and any  $X \subset \mathbb{R}^d$  for  $|X| = n$  finite, there exists an embedding  $f : X \rightarrow \mathbb{R}^m$  for  $m = O(\varepsilon^{-2} \log n)$  such that*

$$\forall x, y \in X, (1 - \varepsilon)\|x - y\|_2^2 \leq \|f(x) - f(y)\|_2^2 \leq (1 + \varepsilon)\|x - y\|_2^2. \quad (5.1)$$

Note one can take the square root of all terms in Eq. (5.1) to say the  $\ell_2$  norm itself is preserved (and not the square), which affects  $\varepsilon$  by roughly a factor of 2; we write the squared version as it makes some later arguments less clumsy.

**Remark 5.0.2.** For two metric spaces  $(X, d_X)$ ,  $(Y, d_Y)$ , a map  $f : X \rightarrow Y$  is often called an *embedding* or *metric embedding*. Let  $\rho$  be the smallest value such that there exists some fixed value  $c$  such that

$$\forall x, y \in X, c \cdot d_Y(f(x), f(y)) \leq d_X(x, y) \leq \rho c \cdot d_Y(f(x), f(y)).$$

If such  $\rho$  exists, then we say it is the *distortion* of  $f$ ; if it does not exist, then we say that  $f$  has unbounded distortion. The JL lemma thus is the statement that any  $n$ -point subset of Euclidean space embeds into  $m$ -dimensional Euclidean space for  $m = O(\varepsilon^{-2} \log n)$  with distortion at most  $1 + \varepsilon$  (specifically as written above, it would be  $\sqrt{(1 + \varepsilon)/(1 - \varepsilon)}$ , but this is  $1 + O(\varepsilon)$  and can be made at most  $1 + \varepsilon$  by changing  $m$  by a constant factor).

A common use of the JL lemma is in the design of approximate algorithms for high-dimensional computational geometry problems. The idea is that given some input  $X$  of a set of high-dimensional vectors, rather than solve the computational problem on  $X$  we can instead solve it on  $f(X)$  for an embedding  $f$  as in the JL lemma. Then presumably, since  $f(X)$  lives in lower dimension, the algorithm is faster.

Consider as one example the  $k$ -means clustering example. In this problem we have input  $X = \{x_1, \dots, x_n\} \subset \mathbb{R}^d$  and a given integer parameter  $k > 1$ , and we would like to return  $y_1, \dots, y_k \in \mathbb{R}^d$  minimizing

$$\sum_{i=1}^n \min_{1 \leq j \leq k} \|x_i - y_j\|_2^2.$$

Any choice of the  $y_j$ 's induce a Voronoi partition  $\mathcal{P} = (\mathcal{P}_1, \dots, \mathcal{P}_k)$  on  $X$ :  $\mathcal{P}_r$  is the set of all  $i$  such

that  $y_r = \operatorname{argmin}_j \|x_i - y_j\|_2^2$ . One can then rewrite the  $k$ -means objective as

$$\min_{k\text{-partitions } \mathcal{P}} \min_{y_1, \dots, y_k} \sum_{i=1}^n \|x_i - y_{\pi_{\mathcal{P}}(i)}\|_2^2 = \min_{k\text{-partitions } \mathcal{P}} \min_{y_1, \dots, y_k} \sum_{j=1}^k \sum_{i \in \mathcal{P}_j} \|x_i - y_j\|_2^2. \quad (5.2)$$

where  $\pi_{\mathcal{P}}(i)$  denotes the partition  $j$  such that  $i \in \mathcal{P}_j$ . One can show then via calculus that for any fixed  $\mathcal{P}$ , the optimal choice of the centers  $y_j$  are the centroids  $\mu_j := (1/|\mathcal{P}_j|) \sum_{i \in \mathcal{P}_j} x_i$ . Thus Eq. (5.2) can be rewritten, and via some simple subsequent manipulations one obtains that the objective is

$$\min_{k\text{-partitions } \mathcal{P}} \sum_{j=1}^k \sum_{i \in \mathcal{P}_j} \|x_i - \mu_j\|_2^2 = \min_{k\text{-partitions } \mathcal{P}} \sum_{j=1}^k \frac{1}{|\mathcal{P}_j|} \sum_{i < i' \in \mathcal{P}_j} \|x_i - x_{i'}\|_2^2. \quad (5.3)$$

From Eq. (5.3) it is apparent that if an embedding  $f$  preserves all squares distances within  $X$  up to  $1 \pm \varepsilon$ , then the cost of *any* clustering (i.e. any partition  $\mathcal{P}$ ) is similarly preserved. Thus we obtain an approximately optimal clustering by mapping  $X$  down to low dimension via the JL lemma, solving  $k$ -means there, then using the discovered clustering partition even for the original high-dimensional  $X$ . Thus, for  $(1 + \varepsilon)$ -approximation it suffices to solve  $k$ -means clustering in dimension  $O(\varepsilon^{-2} \log n)$ .

All known proofs of the JL lemma first prove the following ‘‘Distributional Johnson-Lindenstrauss lemma’’ (DJL lemma).

**Lemma 5.0.3** (DJL lemma). *For any  $\varepsilon, \delta \in (0, 1/2)$  and integer  $d > 1$ , there exists a distribution  $\mathcal{D}_{\varepsilon, \delta}$  over matrices  $\Pi \in \mathbb{R}^{m \times d}$  for  $m = O(\varepsilon^{-2} \log(1/\delta))$  such that for any fixed  $z \in \mathbb{R}^d$  with  $\|z\|_2 = 1$ ,*

$$\mathbb{P}_{\Pi \sim \mathcal{D}_{\varepsilon, \delta}} (|\|\Pi z\|_2^2 - 1| > \varepsilon) < \delta.$$

The JL lemma is then a corollary of the DJL lemma for the following reason: we set  $\delta < 1/n^2$  and pick a random  $\Pi$  as in the DJL lemma. Then for any  $x \neq y \in X$ , we set  $z_{x,y} := (x - y)/\|x - y\|_2$ . The DJL lemma implies  $\mathbb{P}(|\|\Pi z_{x,y}\|_2^2 - 1| > \varepsilon) < \delta$ , which is equivalent to  $\mathbb{P}(|\|\Pi(x - y)\|_2^2 - \|x - y\|_2^2| > \varepsilon \|x - y\|_2^2) < \delta$ . By a union bound, the probability there exists some  $x \neq y \in X$  such that  $\|\Pi(x - y)\|_2^2 \notin [(1 - \varepsilon)\|x - y\|_2^2, (1 + \varepsilon)\|x - y\|_2^2]$  is at most  $\binom{n}{2} \delta < 1$ . Thus there exists a  $\Pi^*$  such that  $\|\Pi^*(x - y)\|_2^2 \in [(1 - \varepsilon)\|x - y\|_2^2, (1 + \varepsilon)\|x - y\|_2^2]$  for all  $x, y \in X$ . We define  $f(x) = \Pi^* x$ .

The main task is thus to prove the DJL lemma.

## 5.1 Proof of the Distributional Johnson-Lindenstrauss lemma

We prove the DJL lemma using the Hanson-Wright inequality, specifically the tail version (see Corollary 1.1.16). We will let  $\mathcal{D}_{\varepsilon, \delta}$  be the distribution over matrices  $\Pi$  with i.i.d. entries  $\Pi_{r,i} = \sigma_{r,i}/\sqrt{m}$ , where the  $\sigma_{r,i}$  are independent Rademachers (i.e. uniform in  $\{-1, 1\}$ ). Define the matrix

$$B_z = \frac{1}{\sqrt{m}} \cdot \begin{bmatrix} -z^\top - & 0 & \dots & 0 \\ 0 & -z^\top - & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & -z^\top - \end{bmatrix}. \quad (5.4)$$

Define the vector  $\sigma \in \{-1, 1\}^{md}$  by  $\sigma = (\sigma_{1,1}, \sigma_{1,2}, \dots, \sigma_{1,d}, \dots, \sigma_{m,1}, \dots, \sigma_{m,d})$ . Then  $\Pi z = B_z \sigma$ . Thus  $\|\Pi z\|_2^2 - 1 = \|B_z \sigma\|_2^2 - 1 = \sigma^\top B_z^\top B_z \sigma - \mathbb{E} \sigma^\top B_z^\top B_z \sigma$ . Defining  $A_z := B_z^\top B_z$  and applying Corollary 1.1.16,

$$\mathbb{P}(|\sigma^\top A_z \sigma - \mathbb{E} \sigma^\top A_z \sigma| > \varepsilon) \lesssim e^{-C\varepsilon^2/\|A_z\|_F^2} + e^{-C\varepsilon/\|A\|}. \quad (5.5)$$

We thus need to bound  $\|A_z\|_F$  and  $\|A_z\|$ .

We see  $A_z$  is an  $md \times md$  block-diagonal matrix with  $m$  blocks, where each block equals  $(1/m)zz^\top$ . The squared Frobenius norm is  $(1/m^2)\sum_{r=1}^m \|zz^\top\|_F^2 = (1/m^2)\sum_{r=1}^m \sum_{i,j} z_i^2 z_j^2 = (1/m)\|z\|_2^4 = 1/m$ . We also have  $\|A_z\|$  is its largest singular value. Since  $A_z$  is real and symmetric, the spectral theorem implies all its eigenvalues are real. Thus the largest singular value is the largest magnitude of any eigenvalue. Since  $A_z$  is block-diagonal, its eigenvalues are the eigenvalues of each block. Thus we just need to bound the largest eigenvalue of  $(1/m)zz^\top$ . This is a rank-1 matrix whose sole eigenvector with nonzero eigenvalue is  $z$ , which has corresponding eigenvalue  $(1/m)\|z\|_2^2 = 1/m$ . Thus overall, [Eq. \(5.5\)](#) is bounded by

$$e^{-C\varepsilon^2/m} + e^{-C\varepsilon/m}$$

which is at most  $\delta$  for  $m = \Omega(\varepsilon^{-1} \log(1/\delta) + \varepsilon^{-2} \log(1/\delta)) = \Omega(\varepsilon^{-2} \log(1/\delta))$ , as desired.

## 5.2 Lower bound

One can ask whether the  $m = \Omega(\varepsilon^{-2} \log n)$  is optimal: does there exist a set  $X$  of size  $n$  such that any  $(1 + \varepsilon)$ -distortion embedding of  $X$  with the Euclidean metric into  $m$ -dimensional Euclidean space must have  $m = \Omega(\varepsilon^{-2} \log n)$ . The first lower bound for this problem was in the original JL paper itself [[JL84](#)], which showed that for  $\varepsilon$  smaller than some fixed constant,  $m = \Omega(\log n)$  is required. We present their argument here since it uses volumetric reasoning that will be useful later for the optimal lower bound in [Subsection 5.2.2](#).

The point set for the original JL lower bound was  $X = \{0, e_1, \dots, e_{n-1}\} \subset \mathbb{R}^n$ , where  $e_i$  is the  $i$ th standard basis vector (1 in the  $i$ th position and 0 elsewhere). Suppose  $f : X \rightarrow \mathbb{R}^m$  satisfies

$$\forall x, y \in X, (1 - \varepsilon)\|x - y\|_2 \leq \|f(x) - f(y)\|_2 \leq (1 + \varepsilon)\|x - y\|_2.$$

Without loss of generality we may assume  $f(0) = 0$  (else translate the image of  $f$  by  $f(0)$ , which does not change its distortion). Define  $\tilde{e}_i := f(e_i)$ . Then since  $\|\tilde{e}_i - 0\|_2 = \|\tilde{e}_i - f(0)\|_2 = (1 \pm \varepsilon)\|e_i\|_2$ , we have that  $1 - \varepsilon \leq \|\tilde{e}_i\|_2 \leq 1 + \varepsilon$ . We also have that for  $i \neq j$ ,  $\|\tilde{e}_i - \tilde{e}_j\|_2 \geq (1 - \varepsilon)\sqrt{2}$ , so the radius  $(1 - \varepsilon)\sqrt{2}/2$  balls around the  $\tilde{e}_i$  have disjoint interior. Meanwhile, each of these balls is fully contained in a radius  $(\max_i \|\tilde{e}_i\|_2) + (1 - \varepsilon)\sqrt{2}/2 \leq (1 + \varepsilon + (1 - \varepsilon)\sqrt{2}/2)$ -radius ball about the origin. Thus, for  $B_d(a, r)$  denoting the ball of radius  $r$  about point  $a$  under metric  $d$ ,

$$\begin{aligned} \text{vol}(B_{\ell_2}(0, 1 + \varepsilon + (1 - \varepsilon)\sqrt{2}/2)) &\geq \text{vol}(\cup_{i=1}^{n-1} B_{\ell_2}(\tilde{e}_i, (1 - \varepsilon)\sqrt{2}/2)) \\ &= \sum_{i=1}^{n-1} \text{vol}(B_{\ell_2}(\tilde{e}_i, (1 - \varepsilon)\sqrt{2}/2)) && \text{(by disjointness of the balls)} \\ &= (n - 1) \cdot \text{vol}(B_{\ell_2}(0, (1 - \varepsilon)\sqrt{2}/2)), \end{aligned}$$

which implies  $n - 1 \leq \left(\frac{1 + \varepsilon + (1 - \varepsilon)\sqrt{2}/2}{(1 - \varepsilon)\sqrt{2}/2}\right)^m$ , so that  $m = \Omega(\log n)$  as long as  $\varepsilon < 1$  is a fixed constant. This follows since the volume ratio of a radius- $r_1$  ball and radius- $r_2$  ball in dimension  $m$  is  $(r_2/r_1)^m$ .

Historically, after the original JL lower bound, [[Alo03](#)] proved an improved lower bound of  $m = \Omega(\min\{n, \varepsilon^{-2} \frac{\log n}{\log(1/\varepsilon)}\})$  that depended on  $\varepsilon$ . It was though slightly suboptimal in  $m$  (by a  $\log(1/\varepsilon)$  factor). The optimal lower bound of  $\Omega(\min\{n, \varepsilon^{-2} \log n\})$  was proven in [[LN16](#)] but only against embeddings  $f$  which are *linear*. The linearity assumption was removed and a lower bound of  $m = \Omega(\varepsilon^{-2} \log n)$  was proven in [[LN17](#)] against all embeddings, even nonlinear ones, under the assumption  $1/\varepsilon^2 \leq \min\{n, d\}$ .<sup>99</sup> This assumption is almost optimal, in the sense that there is

always an upper bound of  $\min\{n-1, d\}$ :  $m = d$  is achievable with  $\varepsilon = 0$  by the identity map, and we can also always translate the points so that one point is zero (this does not change pairwise distances), so that the points span an at most  $(n-1)$ -dimensional subspace, also implying that  $m = n-1$  is achievable with  $\varepsilon = 0$ . Thus we can never hope to prove an  $\varepsilon^{-2} \log n$  lower bound in the case that this expression is larger than  $\min\{n-1, d\}$ . A lower bound of  $m = \Omega(\min\{n, d, \varepsilon^{-2} \log(\varepsilon^2 n)\})$  was then proven in [AK17] for all  $\varepsilon$ .

### 5.2.1 Distributional JL

One can also ask whether the DJL lemma (Lemma 5.0.3) is optimal. Note that even if so, it does not imply that the JL lemma is optimal, because in principle one could achieve improved bounds for Euclidean dimensionality reduction without using random linear maps as in the DJL lemma (though it turns out not be the case, as we will see in Subsection 5.2.2). The DJL lemma was in fact shown optimal in [JW13, KMN11]. The lower bound of [JW13] is via communication complexity and is more general, in that it is a lower bound applying to any communication protocol that allows for estimating  $\ell_2$  norms with small failure probability. We sketch here the lower bound argument of [KMN11].

**Theorem 5.2.1.** *Let  $n, d > 1$  and  $\varepsilon, \delta \in (0, 1/2)$  be fixed. Suppose  $\mathcal{D}_{\varepsilon, \delta}$  is a distribution over  $\mathbb{R}^{m \times d}$  such that for all  $z \in \mathbb{S}^{d-1}$  (the unit sphere),*

$$\mathbb{P}_{\Pi \sim \mathcal{D}_{\varepsilon, \delta}} (|\|\Pi z\|_2^2 - 1| > \varepsilon) < \delta.$$

*Then  $m = \Omega(\min\{d, \varepsilon^{-2} \log(1/\delta)\})$ .*

*Proof sketch.* We use (the easy direction of) Yao's minimax principle. Specifically,

$$\begin{aligned} & \forall z \in \mathbb{S}^{d-1} \quad \mathbb{P}_{\Pi \sim \mathcal{D}_{\varepsilon, \delta}} (|\|\Pi z\|_2^2 - 1| > \varepsilon) < \delta \\ \implies & \mathbb{P}_{z \sim \mathcal{Z}} \mathbb{P}_{\Pi \sim \mathcal{D}_{\varepsilon, \delta}} (|\|\Pi z\|_2^2 - 1| > \varepsilon) < \delta && \text{(for any distribution } \mathcal{Z} \text{ over } \mathbb{S}^{d-1}) \\ \implies & \mathbb{P}_{\Pi \sim \mathcal{D}_{\varepsilon, \delta}} \mathbb{P}_{z \sim \mathcal{Z}} (|\|\Pi z\|_2^2 - 1| > \varepsilon) < \delta \\ \implies & \exists \Pi \in \mathbb{R}^{m \times d} \quad \mathbb{P}_{z \sim \mathcal{Z}} (|\|\Pi z\|_2^2 - 1| > \varepsilon) < \delta. \end{aligned}$$

One then simply shows for an appropriate "hard distribution"  $\mathcal{Z}$ , such  $\Pi$  cannot exist unless  $m = \Omega(\min\{d, \varepsilon^{-2} \log(1/\delta)\})$ . It turns out one can show this for the hard distribution of  $\mathcal{Z}$  simply being a uniform vector over the sphere.  $\square$

### 5.2.2 Optimal JL lower bound

Before we prove the optimal JL lower bound of [LN17], we introduce some necessary background from convex geometry.

**Definition 5.2.2.** A *convex body* is a compact, convex subset of  $\mathbb{R}^d$  with non-empty interior. Recall  $K \subset \mathbb{R}^d$  is *convex* if for all  $x, y \in K$ , the straight line from  $x$  to  $y$  is fully contained in  $K$ . It is *compact* if it is closed (contains all its limit points) and bounded (this is not the general definition of compactness, but it's equivalent for  $\mathbb{R}^d$ ). A convex body is *symmetric* if  $x \in K \Leftrightarrow -x \in K$ . Note a symmetric convex body must thus contain 0 (look at the line between  $x$  and  $-x$  for some  $x \in K$ , which must exist since  $K$  has non-empty interior).

Symmetric convex bodies are in 1-to-1 correspondence with normed spaces in  $\mathbb{R}^d$ . Specifically, given a norm  $\|\cdot\|$  we can define a symmetric convex body  $K = \{x : \|x\| \leq 1\}$ . Given a symmetric convex body  $K$ , we can define a norm  $\|\cdot\|_K$  by  $\|x\|_K = \sup\{t : tx \in K\}$ .

**Definition 5.2.3** (entropy numbers,  $\epsilon$ -nets). Let  $(X, d_X)$  be a metric space and  $T \subseteq X$ . For  $\epsilon \geq 0$ , an  $\epsilon$ -net of  $T$  is a subset  $T' \subseteq X$  such that for all  $x \in T$ , there exists  $x' \in T'$  such that  $d_X(x, x') \leq \epsilon$ . We define the *entropy number*  $\mathcal{N}(T, d_X, \epsilon)$  as the size of the smallest  $\epsilon$ -net of  $T$  under metric  $d_X$ . Given a set  $A$  and symmetric convex body  $K$ , both in  $\mathbb{R}^d$ , we abuse notation and also define  $\mathcal{N}(A, K)$  to be the minimum number of translations of  $K$  required to cover every point in  $A$ . Note that  $\mathcal{N}(T, d_{\|\cdot\|_K}, \epsilon) = \mathcal{N}(T, \epsilon K)$ .

The following lemma is a standard volumetric argument.

**Lemma 5.2.4.** *Let  $B$  be the unit ball of some norm  $\|\cdot\|$  in  $\mathbb{R}^d$ . Then  $\mathcal{N}(B, \|\cdot\|, \epsilon) \leq (1 + 2/\epsilon)^d$ .*

*Proof.* We define the *packing number*  $\mathcal{P}(A, d, \epsilon)$  as the maximum number of radius- $\epsilon$  pairwise disjoint balls under metric that can fit in  $A$ , centered at points in  $A$ . Then consider an optimal packing for  $\mathcal{P}(B, \|\cdot\|, \epsilon/2)$  with centers  $c_1, \dots, c_N$ . The union of all such balls fits inside a ball of radius  $1 + \epsilon/2$ , and thus  $\text{vol}(B_{\|\cdot\|}(0, 1 + \epsilon/2)) \geq N \cdot \text{vol}(B_{\|\cdot\|}(c_i, \epsilon/2))$ . Thus  $N$  is at most the ratio of those volumes, which is  $(1 + 2/\epsilon)^d$ . Meanwhile, the balls of radius  $\epsilon$  about the  $c_i$  must cover  $B$  completely. This is because if there is some  $x \in B$  which is not covered, then  $\forall i \in [N]$ ,  $B_{\|\cdot\|}(x, \epsilon/2) \cap B_{\|\cdot\|}(c_i, \epsilon/2) = \emptyset$ . But then we could have added  $x$  to the packing, violating maximality of the packing, a contradiction.  $\square$

Next, we make a simple observation relating the additive preservation of dot products and Euclidean embeddings with low distortion.

**Lemma 5.2.5.** *Let  $X \subset S^{d-1}$  be such that it contains 0 and  $f : X \rightarrow \mathbb{R}^m$  satisfies  $f(0) = 0$  and*

$$\forall x, y \in X, (1 - \epsilon)\|x - y\|_2^2 \leq \|f(x) - f(y)\|_2^2 \leq (1 + \epsilon)\|x - y\|_2^2. \quad (5.6)$$

*Then*

$$\forall x, y \in X \setminus \{0\}, |\langle f(x), f(y) \rangle - \langle x, y \rangle| \leq 3\epsilon. \quad (5.7)$$

*Proof.* We have

$$\|x - y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 - 2\langle x, y \rangle = 2 - 2\langle x, y \rangle$$

and

$$\|f(x) - f(y)\|_2^2 = \|f(x)\|_2^2 + \|f(y)\|_2^2 - 2\langle f(x), f(y) \rangle = 1 \pm \epsilon + 1 \pm \epsilon - 2\langle f(x), f(y) \rangle.$$

Subtracting and rearranging, we have

$$|\langle f(x), f(y) \rangle - \langle x, y \rangle| \leq \frac{1}{2} [2\epsilon + \epsilon\|x - y\|_2^2] \leq 3\epsilon.$$

The above inequality holds by the triangle inequality since  $\|x - y\|_2 \leq \|x\|_2 + \|y\|_2 = 2$ .  $\square$

We now outline the proof of the lower bound, which is a compression argument (as in [Section 3.1](#)). We define a collection  $\mathcal{X}$  of  $n$ -point sequences in  $\mathbb{R}^d$ . We show that if for each  $X \in \mathcal{X}$  there exists an embedding  $f_X : X \rightarrow \mathbb{R}^m$  satisfying [Eq. \(5.7\)](#) for  $m \ll \epsilon^{-2} \log n$  with error  $\epsilon$ , then there must be an injection  $g : \mathcal{X} \rightarrow \{0, 1\}^S$  for some  $S < \log_2 |\mathcal{X}|$ . This is a contradiction, and thus there must be *some*  $X \in \mathcal{X}$  such that no such  $f_X$  exists. By [Lemma 5.2.5](#), there must thus

be no embedding with distortion at most  $1 + \varepsilon/3$  into such small  $m$ . The theorem then follows by applying the argument to  $\varepsilon' = 3\varepsilon$ .

Define  $k = \lfloor 1/(100\varepsilon^2) \rfloor$  and  $y_S = (1/\sqrt{k}) \sum_{i \in S} e_i$ . We let  $\mathcal{X}$  consist of all point sequences of the form  $X = (0, e_1, e_2, \dots, e_d, y_{S_1}, \dots, y_{S_{n-d-1}})$  where the  $S_i$  are subset of  $[d]$  of size  $k$  each. Thus  $|\mathcal{X}| = \binom{d}{k}^{n-d-1}$ . Note that  $\langle e_i, y_S \rangle$  is either 0 (if  $i \notin S$ ) or  $10\varepsilon$  (if  $i \in S$ ). Thus knowing  $\langle e_i, y_{S_j} \rangle$  for each  $i, j$  is enough to determine  $X$  completely. We will show that if the desired  $f_X$  exists for all  $X \in \mathcal{X}$  (i.e. preserving these dot products up to  $\pm\varepsilon$ ), embedding into dimension  $m$ , then there exists an injection  $g: \mathcal{X} \rightarrow \{0, 1\}^{O(nm)}$ , and thus we must have  $nm = \Omega(\log_2 |\mathcal{X}|) = \Omega(nk \log(d/k))$ . Thus  $m = \Omega(k \log(d/k))$ . We first assume  $d = n/\log(1/\varepsilon)$  and  $1/\varepsilon^2 < n^{.99}$ . Thus we have  $m = \Omega(\varepsilon^{-2} \log n)$ , as desired. We now simply must show how to define such  $g$ .

We would like to set  $g(X)$  to be the concatenation of entries of the sequence of vectors  $f_X(0), f_X(e_1), \dots, f_X(e_d), f_X(y_{S_1}), \dots, f_X(y_{S_{n-d-1}})$ , which would be  $nm$  numbers. Then since  $f_X$  preserves dot products up to additive  $\pm\varepsilon$ , we can know whether  $i \in S_j$  based on whether  $\langle f_X(e_i), f_X(y_{S_j}) \rangle$  is either at most  $\varepsilon$  or at least  $9\varepsilon$ . Unfortunately this does not work, since then  $g$  maps  $X$  to  $nm$  real numbers, but we want a map into  $O(nm)$  bits. We now show how to get around this in progressively better ways. We henceforth refer to  $f_X$  as simply  $f$ .

**Constructing  $g$ : first attempt.** Our first approach is to round each entry of  $f(x)$  to the nearest integer multiple of  $\gamma$  for some small  $\gamma > 0$ . Since  $f(x)$  for  $x \neq 0$  has norm that is  $1 \pm \varepsilon$ , each entry of  $f(x)$  must be in the interval  $[-1 - \varepsilon, 1 + \varepsilon]$ . Thus there are at most  $\lceil (2 + 2\varepsilon)/\gamma \rceil = O(1/\gamma)$  possible multiples of  $\gamma$  that each entry could be rounded to. Thus, after rounding, the the sequence  $f(0), f(e_1), \dots, f(e_d), f(y_{S_1}), \dots, f(y_{S_{n-d-1}})$  consumes at most  $O(nd \log(1/\gamma))$  bits to specify. Meanwhile, the dot products after rounding are

$$\sum_{i=1}^m (f(e_i) \pm \gamma)(f(y_S) \pm \gamma) = \langle f(e_i), f(y_S) \rangle \pm \gamma \|f(e_i)\|_1 \pm \gamma \|f(y_S)\|_1 \pm m\gamma^2$$

By Cauchy-Schwarz,  $\|f(e_i)\|_1 \leq \sqrt{m} \|f(e_i)\|_2 \leq (1 + O(\varepsilon))\sqrt{m}$ , and similarly for  $\|f(y_S)\|_1$ . Thus the total additive error in computing  $\langle f(e_i), f(y_S) \rangle$  is at most  $\varepsilon$  for  $\gamma = \Theta(\varepsilon/\sqrt{m})$ . Recalling  $\langle f(e_i), f(y_S) \rangle$  approximates  $\langle e_i, y_S \rangle$  up to additive  $\pm\varepsilon$ , we thus know  $\langle e_i, y_S \rangle$  up to additive  $\pm 2\varepsilon$ . We can thus declare  $i \in S$  if our estimated dot product is at least  $8\varepsilon$ , and  $i \notin S$  if it is at most  $2\varepsilon$ . Our  $g$  then maps into  $O(nm \log(1/\gamma)) = O(nm \log(m/\varepsilon))$  bits. Since this quantity must be at least  $\Omega(n\varepsilon^{-2} \log n)$ , we have  $m = \Omega(\varepsilon^{-2} \frac{\log n}{\log(1/\varepsilon) + \log \log n})$ .

**Constructing  $g$ : second attempt.** In the first attempt we picked a  $\gamma$ -net  $B'$  of  $B_{\ell_2}(0, 1 + \varepsilon)$  under the  $\ell_\infty$  norm and rounded each  $f(x)$  to the closest element in  $B'$ . What if we instead picked  $B''$  to be a  $\gamma$ -net under  $\ell_2$  and not  $\ell_\infty$ ? Let  $\widetilde{f(x)}$  be the closest point to  $f(x)$  in  $B''$  (under  $\ell_2$ ). Then

$$\begin{aligned} \langle \widetilde{f(x)}, \widetilde{f(y)} \rangle &= \langle f(x) + (\widetilde{f(x)} - f(x)), f(y) + (\widetilde{f(y)} - f(y)) \rangle \\ &= \langle f(x), f(y) \rangle + \langle f(x), \widetilde{f(y)} - f(y) \rangle + \langle \widetilde{f(x)} - f(x), f(y) \rangle + \langle \widetilde{f(x)} - f(x), \widetilde{f(y)} - f(y) \rangle \\ &= \langle f(x), f(y) \rangle + O(\gamma) \end{aligned}$$

with the last inequality holding by Cauchy-Schwarz, since  $\|\widetilde{f(y)} - f(y)\|_2, \|\widetilde{f(x)} - f(x)\|_2 \leq \gamma$ . Thus if we set  $\gamma = c\varepsilon$  for a sufficiently small constant  $c$ , we have  $\langle \widetilde{f(x)}, \widetilde{f(y)} \rangle = \langle f(x), f(y) \rangle \pm \varepsilon$ . Thus again we can thus declare  $i \in S$  if  $\langle \widetilde{f(x)}, \widetilde{f(y)} \rangle$  is at least  $8\varepsilon$ , and  $i \notin S$  if it is at most  $2\varepsilon$ .



By Lemma 5.2.4  $B''$  has size at most  $O(1/\varepsilon)^m$ , and thus the number of bits to specify  $(f(x))_{x \in X}$  is  $O(nm \log(1/\varepsilon))$ , which must be  $\Omega(\log |\mathcal{X}|)$ . Thus  $m = \Omega(\varepsilon^{-2} \frac{\log n}{\log(1/\varepsilon)})$ . This recovers the lower bound of Alon [Alo03], but with a totally different proof.

**Constructing  $g$ : final attempt.** We now remove the  $\log(1/\varepsilon)$  in the denominator in the previous attempt. Recall the second attempt showed that it was enough to know  $\langle \widetilde{f(e_i)}, \widetilde{f(y_{S_j})} \rangle$  for each  $i \in [d], j \in [n - d - 1]$  to recover  $X \in \mathcal{X}$ . Thus if we define  $v_j \in \mathbb{R}^d$  by  $v_j := \langle \widetilde{f(e_i)}, \widetilde{f(y_{S_j})} \rangle$ , then it is enough to know all the  $v_j$ . In fact, it is even enough to know a set of  $\tilde{v}_j$  such that  $\|\tilde{v}_j - v_j\|_\infty \leq \varepsilon$  for each  $j$ . That is, we only need to know roundings of the  $v_j$  in an  $\varepsilon$ -net under the  $\ell_\infty$  metric, since then we can distinguish whether  $i \in S_j$  by whether  $(\tilde{v}_j)_i$  is at most  $3\varepsilon$  or at least  $7\varepsilon$ .

If we define a matrix  $A \in \mathbb{R}^{d \times m}$  to have its  $i$ th row be  $\widetilde{f(e_i)}$ , then we see that  $v_j = Ay_{S_j}$ . Thus the  $v_j$  live in the subspace  $E$  which is the column space of  $A$ , and since  $A$  has at most  $m$  columns, we have  $\dim(E) \leq m$ . Also since  $\langle \widetilde{f(e_i)}, \widetilde{f(y_S)} \rangle = \langle f(e_i), f(y_S) \rangle \pm \varepsilon = \langle e_i, y_S \rangle \pm 2\varepsilon$ , we know that  $\|v_j\|_\infty \leq 12\varepsilon$  for each  $j$ . Thus if we define  $K := E \cap B_{\ell_\infty^d}(0, 12\varepsilon)$ , then  $v_j \in K$  for each  $j$ . Now the beauty:  $K$  is a symmetric convex body, so it defines a norm! Thus if we let  $K'$  be a  $1/12$ th-net of  $K$  in the  $K$ -norm, it implies that for any  $v \in K$  there exists  $\tilde{v} \in K'$  such that  $\|v - \tilde{v}\|_\infty \leq \varepsilon$ . This will be our definition of the  $\tilde{v}_j$ : simply the rounding of the  $v_j$  to the closest points (under the  $K$ -norm) in  $K'$ . The size of  $K'$  is  $O(1 + 2/(1/12))^m = O(1)^m$  since  $K$  is  $m$ -dimensional (Lemma 5.2.4), and thus  $\tilde{v}_j$  can be specified in  $O(m)$  bits. Thus we all  $\tilde{v}_j$  for every  $j \in [n - d - 1]$  can be specified in  $O(nm)$  bits combined!

There is one slight catch: recall in compression arguments we typically show that the “compression” is an injection by showing that it is possible to invert. But for the decompressor to invert, they will need to know which body  $K$  we are talking about. For that, they need to know the matrix  $A$ , which depends on  $f = f_X$  and thus depends on  $X$ . We accomplish this by simply writing down  $A$  explicitly, row-by-row; each  $\widetilde{f(e_i)}$  takes  $O(m \log(1/\varepsilon))$  bits, so in total this takes  $O(dm \log(1/\varepsilon)) = O(nm)$  bits since  $d = n/\log(1/\varepsilon)$ .

**Removing the assumption  $d = n/\log(1/\varepsilon)$ .** Above we showed how to obtain the optimal lower bound when  $d = n/\log(1/\varepsilon)$ . What about for other  $d$ ? Showing a hard set with larger  $d$  exists is easy: we simply take the hard point set in dimension  $n/\log(1/\varepsilon)$  then pad all the vectors with zeroes to make them dimension  $d$ .

What about for smaller  $d$ ? Suppose  $X$  is a hard point set in dimension  $n/\log(1/\varepsilon)$ , i.e. any embedding with distortion  $(1 + \varepsilon)^2 = 1 + O(\varepsilon)$  needs target dimension  $\geq c\varepsilon^{-2} \log n$ . Consider the set of points  $f(X)$  in dimension  $d' = O(\varepsilon^{-2} \log n)$  where  $f: X \rightarrow \mathbb{R}^{d'}$  has distortion at most  $1 + \varepsilon$ ; such  $f$  exists by the JL lemma. Then  $f(X)$  must similarly be “hard” (i.e. require target dimension  $\geq c\varepsilon^{-2} \log n$  to achieve distortion  $1 + \varepsilon$ ), since otherwise if one could embed into such small dimension with low distortion via a mapping  $g$ , then  $g \circ f$  would be a good embedding for  $X$ .

### 5.3 Speeding up Johnson-Lindenstrauss transforms

Recall one motivation for dimensionality reduction: we have some dataset  $X \subset \mathbb{R}^d$  for  $d$  large on which we would like to solve some computational geometry problem. We would like to map  $X$  down to some  $f(X) \subset \mathbb{R}^m$  for  $m \ll d$  then solve the problem on  $f(X)$ , which is hopefully faster than solving it on  $X$  since  $m$  is much smaller. Achieving  $m$  as small as possible is thus important to solving  $f(X)$  quickly, but also important is being able to compute  $f(X)$  quickly. That is, given a point  $x \in X$ , we would like an  $f$  that allows us to compute  $f(x)$  quickly.

The original JL map of [JL84] chose  $f(z) = \Pi z$ , where  $\Pi$  is (scaled) orthogonal projection onto a random  $m$ -dimensional subspace of  $\mathbb{R}^d$ . Thus we would like to pick a random basis of  $m$  orthonormal vectors. We can accomplish this by picking a gaussian vector  $g_1 \sim \mathcal{N}(0, I_d)$  then letting the first row of  $\Pi$  be  $r_1 := g_1/\|g_1\|_2$ . For  $r_2$ , we pick another  $g_2 \sim \mathcal{N}(0, I_d)$  independently and define  $\tilde{g}_2 = g_2 - \langle g_2, g_1 \rangle g_1$  then  $r_2 := \tilde{g}_2/\|\tilde{g}_2\|_2$ . That is, we first subtract its projection onto  $g_1$  before normalizing to form  $g_2$ . In general for  $r_j$  we pick  $g_j$ , subtract off its projection onto the span of  $g_1, \dots, g_{j-1}$ , then normalize (the ‘‘Gram-Schmidt process’’). In the end,  $\Pi$  is a dense, unstructured matrix, so computing  $\Pi z$  takes  $O(md)$  time. Can we do better?

The first work to provide any improvement in embedding time is Achlioptas [Ach03]. That work looked at the bound  $m = C\varepsilon^{-2} \ln(1/\delta)$  for DJL and provided a new argument showing that the best known bound of  $C$  is achievable while simultaneously making the matrix  $\Pi$  sparser by a factor of three. That is, the  $\Pi_{i,j}$  are independent, equal to  $1/\sqrt{s}$  with probability  $1/6$ ,  $-1/\sqrt{s}$  with probability  $1/6$ , and 0 with probability  $2/3$ , for  $s = m/3$ . Thus in expectation,  $2/3$  of the entries of  $\Pi$  are non-zero. Since  $\Pi z = \sum_i z_i \Pi^i$  where  $\Pi^i$  is the  $i$ th column of  $\Pi$ , one could potentially perform this computation faster than if  $\Pi$  had been completely dense.

The next improvement came from the ‘‘Fast Johnson-Lindenstrauss Transform’’ (FJLT) of [AC09]. The main idea here is that if  $S \in \mathbb{R}^{m \times d}$  is a sampling matrix, i.e. a 1 in a random location per row (and zeroes elsewhere in the row), with rows chosen independently, then  $\mathbb{E} \|\frac{1}{\sqrt{m}} S z\|_2^2 = \|z\|_2^2$  for any vector  $z$ . Note computing the mapping  $z \mapsto \frac{1}{\sqrt{m}} S z$  is fast ( $O(m)$  time). The problem is that the variance might be quite high, e.g. if  $z$  has its mass concentrated on one or few coordinates. The FJLT applies a random pre-conditioning operation to  $z$ , i.e.  $z \mapsto R z$  for a certain random orthogonal matrix  $R$ , such that  $\|R z\|_\infty / \|R z\|_2$  is small with high probability, which is one mathematical way to express that  $R z$  is ‘‘well-spread’’, with no one coordinate having too much mass. One can then show that conditioning on the event that it is well-spread,  $\frac{1}{\sqrt{m}} S R z$  has roughly the same norm as  $z$  with high probability. The runtime of the FJLT to embed a vector  $z$  is  $O(d \log d + m^3)$ , though subsequent works have improved the  $m^3$  term [AL09, AL13, KW11]. We discuss the FJLT further in Subsection 5.3.2.

The main downside of the FJLT is that it is not fast for sparse inputs  $z$ . Sparse data comes up frequently in data science applications, e.g. representing a document as a ‘‘bag of words’’ (that is,  $d$  may be the size of the dictionary, and for some document  $D$ ,  $z_D \in \mathbb{R}^d$  is defined by  $z_i$  being the number of occurrences of word  $i$  in  $D$ ). Since most documents do not contain the entire dictionary,  $z$  is likely sparse. The FJLT works by preconditioning  $z$  to eliminate its sparsity with high probability, so that sampling works. But if  $z$  was sparse to begin with, we may hope for running times that depend on its sparsity as opposed to its dimensionality. One way to achieve this is to let the embedding matrix  $\Pi$  itself be sparse. If  $\Pi$  has at most  $s$  non-zero entries per column, then  $\Pi z$  can be computed in time  $O(s \cdot \|z\|_0)$ , where  $\|z\|_0 := |\{i : z_i \neq 0\}|$ . It turns out that for the DJL lemma,  $m = O(\varepsilon^{-2} \log(1/\delta))$  can still be achieved with  $s$  only  $O(\varepsilon m) = O(\varepsilon^{-1} \log(1/\delta))$ . Thus using that  $\Pi z = \sum_i z_i \Pi^i$ , we can compute  $\Pi z$  faster by a factor  $1/\varepsilon$  compared with dense  $\Pi$ . The first work providing an asymptotically sparser embedding matrix was [DKS10], with improvements later in [KN10, BOR10]. The SJLT bound stated above, which used a different construction from the prior works, is from [KN14], with a nearly matching lower bound shown in [NN13b]. A simpler proof of the upper bound was later provided in [CJN18]. We discuss the SJLT further in Subsection 5.3.1.

### 5.3.1 Sparse Johnson-Lindenstrauss Transform

As mentioned, one natural way to speed up JL is to make  $\Pi$  sparse. If  $\Pi$  has  $s$  non-zero entries per column, then  $\Pi x$  can be multiplied in time  $O(s \cdot \|x\|_0)$ . The goal is then to make  $s, m$  as small as

possible.

We consider the `CountSketch` matrix introduced in [Subsection 4.1.2](#), which was first shown to provide DJL in [\[TZ12\]](#) as mentioned in [Remark 4.3.1](#). In the construction  $\Pi$  being the `CountSketch`, one picks a hash function  $h : [d] \rightarrow [m]$  from a 2-wise independent family, and a sequence of bits  $\sigma_1, \dots, \sigma_d \in \{-1, 1\}$  from a 4-wise independent family. Then for each  $i \in [d]$ ,  $\Pi_{h(i),i} = \sigma_i$ , and the rest of the  $i$ th column is 0. It was shown in [\[TZ12\]](#) that this distribution provides DJL for  $m \gtrsim 1/(\varepsilon^2\delta)$ . Note that the column sparsity is  $s = 1$  as described here. The analysis is simply via Chebyshev's inequality, after doing an expectation and variance calculation.

The reason for the poor dependence in  $m$  on the failure probability  $\delta$  is that we use Chebyshev's inequality. We will make improvements by using Hanson-Wright, i.e. a bound on the  $p$ -norms of quadratic forms. Recall that a bound on  $p$ -norms gives tail bounds via Markov's inequality.

To improve the dependence of  $m$  on  $1/\delta$ , we allow ourselves to increase  $s$ . Here we analyze the Sparse JL Transform (SJLT) [\[KN14\]](#). This is a JL distribution over  $\Pi$  having exactly  $s$  non-zero entries per column. The analysis we give below can be found in [\[CJN18\]](#).

As previously, without loss of generality we assume  $z \in \mathbb{R}^d$  has  $\|z\|_2 = 1$ . Our random  $\Pi \in \mathbb{R}^{m \times d}$  satisfies  $\Pi_{r,i} = \eta_{r,i}\sigma_{r,i}/\sqrt{s}$  for some integer  $1 \leq s \leq m$ . The  $\sigma_{r,i}$  are independent Rademachers. The  $\eta_{r,i}$  are Bernoulli random variables satisfying:

- For all  $r, i$ ,  $\mathbb{E} \eta_{r,i} = s/m$ .
- For any  $i$ ,  $\sum_{r=1}^m \eta_{r,i} = s$ . That is, each column of  $\Pi$  has *exactly*  $s$  non-zero entries.
- The  $\eta_{r,i}$  are negatively correlated. That is, for any subset  $S$  of  $[m] \times [n]$ , we have  $\mathbb{E} \prod_{(r,i) \in S} \eta_{r,i} \leq \prod_{(r,i) \in S} \mathbb{E} \eta_{r,i} = (s/m)^{|S|}$ .

The above is satisfied by the `CountSketch`, but any distribution satisfying the above criteria would do. For example, we could also select exactly  $s$  entries per column to be non-zero, uniformly without replacement.

We would like to show the following, which is the main theorem of [\[KN14\]](#).

**Theorem 5.3.1.** *As long as  $m \simeq \varepsilon^{-2} \log(1/\delta)$  and  $s \simeq \varepsilon m$ ,*

$$\forall z : \|z\|_2 = 1, \mathbb{P}_{\Pi}(\|\Pi z\|_2^2 - 1 > \varepsilon) < \delta. \quad (5.8)$$

*Proof.* Abusing notation and treating  $\sigma$  as an  $md$ -dimensional vector,

$$E = \|\Pi z\|_2^2 - 1 = \frac{1}{s} \sum_{r=1}^m \sum_{i \neq j} \eta_{r,i} \eta_{r,j} \sigma_{r,i} \sigma_{r,j} z_i z_j \stackrel{\text{def}}{=} \sigma^\top A_{z,\eta} \sigma,$$

Thus by Hanson-Wright

$$\|E\|_p \leq \sqrt{p} \cdot \|A_{z,\eta}\|_F + p \cdot \|A_{z,\eta}\|_p \leq \sqrt{p} \cdot \| \|A_{z,\eta}\|_F \|_p + p \cdot \| \|A_{z,\eta}\| \|_p.$$

$A_{z,\eta}$  is a block diagonal matrix with  $m$  blocks, where the  $r$ th block is  $(1/s)z^{(r)}(z^{(r)})^\top$  but with the diagonal zeroed out. Here  $z^{(r)}$  is the vector with  $(z^{(r)})_i = \eta_{r,i}z_i$ . Now we just need to bound  $\| \|A_{z,\eta}\|_F \|_p$  and  $\| \|A_{z,\eta}\| \|_p$ , where here the  $p$ -norm is over the randomness in  $\eta$ .

Since  $A_{z,\eta}$  is block-diagonal, its operator norm is the largest operator norm of any block. The eigenvalue of the  $r$ th block is at most  $(1/s) \cdot \max\{\|z^{(r)}\|_2^2, \|z^{(r)}\|_\infty^2\} \leq 1/s$ , and thus  $\|A_{z,\eta}\| \leq 1/s$  with probability 1.

Next, define  $Q_{i,j} = \sum_{r=1}^m \eta_{r,i} \eta_{r,j}$  so that

$$\|A_{z,\eta}\|_F^2 = \frac{1}{s^2} \sum_{i \neq j} z_i^2 z_j^2 \cdot Q_{i,j}.$$

We will show for  $p \simeq s^2/m$  that for all  $i, j$ ,  $\|Q_{i,j}\|_p \lesssim p$ , where we take the  $p$ -norm over  $\eta$ . Therefore for this  $p$ ,

$$\begin{aligned} \|\|A_{z,\eta}\|_F\|_p &= \|\|A_{z,\eta}\|_F^2\|_p^{1/2} \\ &\leq \left\| \frac{1}{s^2} \sum_{i \neq j} z_i^2 z_j^2 \cdot Q_{i,j} \right\|_p^{1/2} \\ &\leq \frac{1}{s} \left( \sum_{i \neq j} z_i^2 z_j^2 \cdot \|Q_{i,j}\|_p \right)^{1/2} && \text{(triangle inequality)} \\ &\leq \frac{1}{\sqrt{m}} \end{aligned}$$

Then by Markov's inequality and the settings of  $p, s, m$ ,

$$\mathbb{P}(\|\| \Pi x \|_2^2 - 1\| > \varepsilon) = \mathbb{P}(|\sigma^\top A_{z,\eta} \sigma| > \varepsilon) < \varepsilon^{-p} \cdot C^p (m^{-p/2} + s^{-p}) < \delta.$$

We now show  $\|Q_{i,j}\|_p \lesssim p$ , for which we use Bernstein's inequality ([Theorem 1.1.17](#)).

Suppose  $\eta_{a_1,i}, \dots, \eta_{a_s,i}$  are all 1, where  $a_1 < a_2 < \dots < a_s$ . Now, note  $Q_{i,j}$  can be written as  $\sum_{t=1}^s Y_t$ , where  $Y_t$  is an indicator random variable for the event that  $\eta_{a_t,j} = 1$ . The  $Y_t$  are not independent, but for any integer  $p \geq 1$  their  $p$ th moment is upper bounded by the case that the  $Y_t$  are independent Bernoulli each of expectation  $s/m$  (this can be seen by simply expanding  $(\sum_t Y_t)^p$  then comparing with the independent Bernoulli case monomial by monomial in the expansion). Thus Bernstein applies, and as desired we have

$$\|Q_{i,j}\|_p = \left\| \sum_t Y_t \right\|_p \lesssim \sqrt{s^2/m} \cdot \sqrt{p} + p \simeq p.$$

□

### 5.3.2 Fast Johnson-Lindenstrauss Transform

Another approach for obtaining fast JL was investigated by Ailon and Chazelle [[AC09](#)]. This approach gives a running time to compute  $\Pi x$  of roughly  $O(d \log d)$ , which is faster than the sparse JL approach when  $x$  is sufficiently dense. They called their transformation the *Fast Johnson-Lindenstrauss Transform (FJLT)*. A construction similar to theirs, which we will analyze here, is the  $m \times d$  matrix  $\Pi$  defined as

$$\Pi = \frac{1}{\sqrt{m}} SHD \tag{5.9}$$

where  $S$  is an  $m \times d$  sampling matrix with replacement (each row has a 1 in a uniformly random location and zeroes elsewhere, and the rows are independent),  $H$  is an *unnormalized bounded orthonormal system*, and  $D = \text{diag}(\alpha)$  for a vector  $\alpha$  of  $n$  independent Rademachers. An unnormalized bounded orthonormal system is a matrix  $H \in \mathbb{R}^{n \times n}$  such that  $H^\top H = d \cdot I$  and  $\max_{i,j} |H_{i,j}| \leq 1$ . For example,  $H$  can be the unnormalized Fourier matrix or Hadamard matrix<sup>1</sup>. The original FJLT

<sup>1</sup>The  $d \times d$  unnormalized Hadamard matrix for  $d = 2^k$  when using 0-based indexing for entries has  $H_{i,j} = (-1)^{\langle \vec{i}, \vec{j} \rangle \bmod 2}$ , where  $\vec{i}$  is the  $k$ -dimensional binary vector obtained by writing  $i$  in binary,  $0 \leq i < d$ .

replaced  $S$  with a random sparse matrix  $P$ , which has certain advantages; see [Remark 5.3.3](#). For a vector  $v \in \mathbb{R}^d$ ,  $\text{diag}(v)$  is a  $d \times d$  diagonal matrix with  $i$ th diagonal entry equal to  $v_i$ .

The motivation for the construction [Eq. \(5.9\)](#) is speed:  $D$  can be applied in  $O(d)$  time,  $H$  in  $O(d \log d)$  time (e.g. using the Fast Fourier Transform for the DFT, or a divide-and-conquer algorithm for the Hadamard transform), and  $S$  in  $O(m)$  time. Thus, overall, applying  $\Pi$  to any fixed vector  $x$  takes  $O(d \log d)$  time. Compare this with using a dense matrix of Rademachers, which takes  $O(md)$  time to apply.

The intuition for why the construction works is as follows. We would ideally like to simply set  $\Pi = \frac{1}{\sqrt{m}}S$  since then  $\mathbb{E} \|\Pi z\|_2^2 = \|z\|_2^2$ . The problem is that the variance may be very large, which would necessitate us setting  $m$  to be large — for example, consider  $z = e_1$ , or more generally any vector with most of its  $\ell_2$  mass concentrated in one or few coordinates. Unless  $m = \Omega(d)$ , it is likely that  $Sz = 0$ . To fix this, [\[AC09\]](#) performs a randomized pre-conditioning step motivated by the so-called “uncertainty principle” from quantum mechanics: a vector and its Fourier transform cannot both be concentrated in few coordinates. Thus rather than apply  $S$  to  $z$ , we could perhaps apply it to  $Hz$  where  $H$  is a bounded orthonormal system (like the Discrete Fourier Transform). There is still a problem though: we may have the opposite problem, namely that  $z$  has its mass well spread across all its coordinates, but  $Hz$  has its mass all collapsed to one or few coordinates. In fact, for the Hadamard transform it is possible to construct example  $z$  in which both  $z$  and  $Hz$  have their mass equally spread on  $\sqrt{d}$  coordinates, so that either one would require  $m = \Omega(\sqrt{d})$  if using a sampling matrix. The work [\[AC09\]](#) fixes this by introducing randomness into the pre-conditioning step, namely by first multiplying by the matrix  $D = \text{diag}(\alpha)$  above. They then show that with high probability  $\|HDz\|_\infty$  is small; conditioned on this event,  $S$  then preserves the norm of  $HDz$  with high probability. We show below the proof of [\[AC09\]](#) that for  $m \gtrsim \varepsilon^{-2} \log(1/\delta) \log(d/\delta)$ , the random  $\Pi$  described in [Eq. \(5.9\)](#) provides DJL.

**Theorem 5.3.2.** *Let  $z \in \mathbb{R}^d$  be an arbitrary unit norm vector, and suppose  $0 < \varepsilon, \delta < 1/2$ . Also let  $\Pi = SHD$  as described above with a number of rows equal to  $m \gtrsim \varepsilon^{-2} \log(1/\delta) \log(d/\delta)$ . Then*

$$\mathbb{P}_{\Pi}(|\|\Pi x\|_2^2 - 1| > \varepsilon) < \delta.$$

*Proof.* Define  $y = HDz$ . Define the event  $\mathcal{E}$  that  $\|y\|_\infty \leq \sqrt{2 \ln(4d/\delta)}$ . Note  $y_i = \sum_{j=1}^d H_{i,j} \alpha_j z_j$ . Then by Khintchine’s inequality ([Theorem 1.1.7](#)),

$$\mathbb{P}_{\alpha}(|y_i| > \sqrt{2 \ln(4d/\delta)}) < 2e^{-\frac{2 \ln(4d/\delta)}{2(\sum_{j=1}^d |H_{i,j}|^2 z_j^2)}} = \frac{\delta}{2d}$$

since  $|H_{i,j}| = 1$ . Thus by a union bound over all  $i \in [d]$ ,  $\mathbb{P}(\mathcal{E}) \geq 1 - \delta/2$ . We now upper bound the conditional probability  $\mathbb{P}(|\|\frac{1}{\sqrt{m}}Sy\|_2^2 - 1| > \varepsilon \mid \mathcal{E})$ , for which we use Bernstein’s inequality ([Corollary 1.1.18](#)).

In the language of [Corollary 1.1.18](#),  $\|\frac{1}{\sqrt{m}}Sy\|_2^2$  can be expressed as  $X := \sum_{i=1}^m X_i$  where the  $X_i$  are i.i.d., each equal to  $y_i^2/m$  for a uniformly random  $i \in [d]$ , and bounded by  $K := 2 \ln(4d/\delta)/m$  when we condition on  $\mathcal{E}$ . Since  $\|y\|_2^2 = d$ , we have  $\mathbb{E} X_i = 1/m$  so that  $\mathbb{E} X = 1$ . We also have

$$\begin{aligned} \mathbb{E} X^2 &= \sum_i \mathbb{E} X_i^2 + \sum_{i \neq j} (\mathbb{E} X_i)(\mathbb{E} X_j) \\ &= m \cdot \mathbb{E} X_1^2 + m(m-1)(\mathbb{E} X_1)^2 \\ &\leq K + 1 \end{aligned}$$

Since  $\sigma^2 := \text{Var}[X] = \mathbb{E}X^2 - (\mathbb{E}X)^2$ , we thus have  $\sigma^2 \leq K$ . Therefore by [Corollary 1.1.18](#),

$$\mathbb{P}\left(\left|\left\|\frac{1}{\sqrt{m}}Sy\right\|_2^2 - 1\right| > \varepsilon \mid \mathcal{E}\right) \lesssim e^{-C\frac{\varepsilon^2}{K}} + e^{-C\frac{\varepsilon}{K}} = e^{-C\frac{\varepsilon^2 m}{2\ln(4d/\delta)}} + e^{-C\frac{\varepsilon m}{2\ln(4d/\delta)}},$$

which is less than  $\delta/2$  for  $m \gtrsim \varepsilon^{-2} \log(1/\delta) \log(d/\delta)$ . Therefore

$$\mathbb{P}\left(\left|\left\|\frac{1}{\sqrt{m}}SHDz\right\|_2^2 - 1\right| \leq \varepsilon\right) \geq \mathbb{P}(\mathcal{E}) \cdot \mathbb{P}\left(\left|\left\|\frac{1}{\sqrt{m}}SHDz\right\|_2^2 - 1\right| \leq \varepsilon \mid \mathcal{E}\right) \geq (1 - \delta/2)^2 > 1 - \delta.$$

□

**Remark 5.3.3.** Note that the FJLT as analyzed above provides suboptimal  $m$ . If one desired optimal  $m$ , one can instead use the embedding matrix  $\Pi'\Pi$ , where  $\Pi$  is the FJLT and  $\Pi'$  is, say, a dense matrix with Rademacher entries having the optimal  $m' = O(\varepsilon^{-2} \log(1/\delta))$  rows. The downside is that the runtime to apply our embedding worsens by an additive  $m \cdot m'$ . [\[AC09\]](#) slightly improved this additive term (by an  $\varepsilon^2$  multiplicative factor) by replacing the matrix  $S$  with a random sparse matrix  $P$ .

**Remark 5.3.4.** It is possible to improve the analysis to show that in fact  $m \gtrsim \varepsilon^{-2} \log(1/\delta) \log(1/(\varepsilon\delta))$  suffices; see Theorem 9 of the full version of [\[CNW16\]](#).

### 5.3.3 Krahmer-Ward theorem

Following work of Ailon and Liberty [\[AL13\]](#), Krahmer and Ward [\[KW11\]](#) proved a theorem that shows that there is a totally different path to analyzing JL matrices of the form  $\frac{1}{\sqrt{m}}SHD$  as in [Subsection 5.3.2](#). The analysis of [\[AC09\]](#) can be summarized as: for fixed  $z$ , condition on the event  $\mathcal{E}$  that  $HD$  is “nice”, then show that sampling preserves the norm with high probability given  $\mathcal{E}$  (here “nice” meant that  $\|HDz\|_\infty$  was small). Krahmer and Ward parenthesize the conditioning in a different way: they instead show that if  $\frac{1}{\sqrt{m}}SH$  is “nice”, then  $\frac{\sqrt{m}}{S}HD$  provides a JL distribution  $\mathcal{D}_{\varepsilon,\delta}$  as in [Lemma 5.0.3](#). The notion of niceness here depends on a notion known as the *restricted isometry property* [\[CT05\]](#), which we elaborate on in a later chapter when we discuss compressed sensing.

**Definition 5.3.5.** A matrix  $\Pi \in \mathbb{R}^{m \times d}$  satisfies the  $(\varepsilon, k)$ -*Restricted Isometry Property* (or  $(\varepsilon, k)$ -*RIP*) if for all  $k$ -sparse vectors  $z \in \mathbb{R}^d$  (i.e.  $|\{i : z_i \neq 0\}| \leq k$ ),

$$(1 - \varepsilon)\|z\|_2^2 \leq \|\Pi z\|_2^2 \leq (1 + \varepsilon)\|z\|_2^2.$$

**Lemma 5.3.6.** For  $\Pi \in \mathbb{R}^{m \times d}$  and  $i \in [d]$ , let  $\Pi^i$  denote the  $m \times d$  matrix where all columns except the  $i$ th one are zeroed out. For  $S \subseteq [d]$  let  $\Pi^S$  denote  $\sum_{i \in S} \Pi^i$ . Then if  $\Pi$  satisfies  $(\varepsilon, k)$ -RIP,

$$(a) \sup_{|S| \leq k} \|(\Pi^S)^\top \Pi^S - I^S\| \leq \varepsilon$$

$$(b) \sup_{\substack{|S|+|T| \leq k \\ S \cap T = \emptyset}} \|(\Pi^S)^\top \Pi^T\| \leq \varepsilon$$

*Proof.* **(a).** Since  $\Pi z = \Pi^{\text{support}(z)} z$ , this holds since for real symmetric  $M$  we have  $\|M\| = \sup_{\|z\|_2=1} |z^\top M z|$ .

**(b).** Fix  $S, T$  with  $|S| + |T| \leq k$ . Then  $\|(\Pi^S)^\top \Pi^T\| = \sup_{\substack{\text{support}(z) \subset S, \|z\|_2=1 \\ \text{support}(w) \subset T, \|w\|_2=1}} z^\top \Pi^\top \Pi w$ . Note

$$z^\top \Pi^\top \Pi w = \langle \Pi z, \Pi w \rangle$$

$$\begin{aligned}
&= \frac{1}{4} (\|\Pi(z+w)\|_2^2 - \|\Pi(z-w)\|_2^2) \\
&= \frac{1}{4} ((1 \pm \varepsilon)\|z+w\|_2^2 - (1 \pm \varepsilon)\|z-w\|_2^2) && (\Pi \text{ satisfies } (\varepsilon, k)\text{-RIP}) \\
&= \frac{1}{4} ((1 \pm \varepsilon)(\|z\|_2^2 + \|w\|_2^2 - \langle z, w \rangle) - (1 \pm \varepsilon)(\|z\|_2^2 + \|w\|_2^2 - 2\langle z, w \rangle)) \\
&= \frac{1}{4} (\pm 2\varepsilon\|z\|_2^2 \pm 2\varepsilon\|w\|_2^2) && (\langle z, w \rangle = 0) \\
&= \pm \varepsilon.
\end{aligned}$$

□

The following is a known result of Haviv and Regev [HR16], following [CT06, RV08, Bou14], which we will not cover the proof of here.

**Theorem 5.3.7.** *Let  $H \in \mathbb{R}^{d \times d}$  be a bounded orthonormal system, and let  $S \in \mathbb{R}^{m \times d}$  be a sampling matrix: its rows are independent, with each row having a 1 in a uniformly random location and zeroes elsewhere. Then for  $m = O(\varepsilon^{-2} \log^2(1/\varepsilon) k \log^2(k/\varepsilon) \log d)$ ,  $\frac{1}{\sqrt{m}}SH$  satisfies  $(\varepsilon, k)$ -RIP with probability at least  $1 - \exp(-\Omega(\log(k/\varepsilon) \cdot \log d))$ .*

**Theorem 5.3.8.** *The following holds for some universal constants  $c, C > 0$ . Suppose  $\varepsilon, \delta \in (0, 1/2)$  are such that  $\Pi \in \mathbb{R}^{m \times d}$  satisfies  $(c\varepsilon, 2k)$ -RIP for  $k = \log(4/\delta)$ . For vector  $a \in \mathbb{R}^d$  let  $D_a$  denote the  $d \times d$  diagonal matrix with  $D_{i,i} = a_i$ . Let  $\sigma \in \{-1, 1\}^d$  be chosen uniformly at random. Then*

$$\forall \|z\|_2 = 1, \mathbb{P}(\left| \|\Pi D_\sigma z\|_2^2 - 1 \right| > \varepsilon) < \delta.$$

*Proof.* Relabel coordinates so that  $|z_1| \geq |z_2| \geq \dots \geq |z_d|$ . Define  $S_i = \{(k-1) \cdot i + 1, (k-1) \cdot i + 2, \dots, (k-1) \cdot i + k\}$  and define  $z_{(i)} := z_{S_i} \in \mathbb{R}^d$ . That is, we partition the entries of  $z$  into size- $k$  blocks, so  $z_{(1)}$  contains the largest  $k$  entries in magnitude,  $z_{(2)}$  contains the next largest  $k$ , etc., with all other entries zeroed out. Also define  $M_{(i)} := M^{S_i}$  for a matrix  $i$  (using the definition of  $M^{S_i}$  from Lemma 5.3.6). Then observing  $D_\sigma z = D_z \sigma$ , we have

$$\begin{aligned}
\|\Pi D_\sigma z\|_2^2 &= \|\Pi D_z \sigma\|_2^2 \\
&= \sigma^\top D_z \Pi^\top \Pi D_z \sigma \\
&= \sigma^\top X \sigma
\end{aligned}$$

for  $X := D_z \Pi^\top \Pi D_z$ , which is a  $d \times d$  matrix. Note  $D_z = \sum_i D_{z(i)}$  and  $\Pi = \sum_i \Pi_{(i)}$ , so that  $D_z \Pi^\top \Pi D_z = \sum_{i,j} D_{z(i)} \Pi_{(i)}^\top \Pi_{(j)} D_{z(j)}$ . We decompose  $X = A + B + B^\top + C$  where

$$\begin{aligned}
A &:= \sum_i D_{z(i)} \Pi_{(i)}^\top \Pi_{(i)} D_{z(i)} \\
B &:= \sum_j D_{z(1)} \Pi_{(1)}^\top \Pi_{(j)} D_{z(j)} \\
C &:= \sum_{\substack{i,j>1 \\ i \neq j}} D_{z(i)} \Pi_{(i)}^\top \Pi_{(j)} D_{z(j)}
\end{aligned}$$

Graphically we have the following picture, where each matrix is filled with zeroes in the white regions, and the portions with black ground have entries identical to that of  $X$ :

$A =$ 
 $B =$ 
 $B^\top =$ 
 $C =$

We will show that for each one of the following events, it fails to hold with probability at most  $\delta/4$ ; the theorem then follows by a union bound:

$$\mathcal{E}_1: \sigma^\top A \sigma = 1 \pm \frac{\varepsilon}{4}$$

$$\mathcal{E}_2: \sigma^\top B \sigma = \pm \frac{\varepsilon}{4}$$

$$\mathcal{E}_3: \sigma^\top B^\top \sigma = \pm \frac{\varepsilon}{4}$$

$$\mathcal{E}_4: \sigma^\top C \sigma = \pm \frac{\varepsilon}{4}$$

**Analyzing  $\mathcal{E}_1$ :**

$$\begin{aligned}
\sigma^\top A \sigma &= \sum_i \sigma_{(i)}^\top D_{z(i)} \Pi_{(i)}^\top \Pi_{(i)} D_{z(i)} \sigma_{(i)} \\
&= \sum_i z_{(i)}^\top D_{\sigma(i)} \Pi_{(i)}^\top \Pi_{(i)} D_{\sigma(i)} z_{(i)} \\
&= \sum_i \|\Pi_{(i)} D_{\sigma(i)} z_{(i)}\|_2^2 \\
&= \left(1 \pm \frac{\varepsilon}{4}\right) \sum_i \|D_{\sigma(i)} z_{(i)}\|_2^2 && (D_{\sigma(i)} z_{(i)} \text{ is } k\text{-sparse}) \\
&= \left(1 \pm \frac{\varepsilon}{4}\right) \sum_i \|z_{(i)}\|_2^2 \\
&= \left(1 \pm \frac{\varepsilon}{4}\right) \|z\|_2^2 \\
&= 1 \pm \frac{\varepsilon}{4}.
\end{aligned}$$

Thus  $\mathcal{E}_1$  holds with probability 1 (given that  $\Pi$  is RIP).

**Analyzing  $\mathcal{E}_2$ :** Let  $x_{(-1)}$  denote  $x - x_{(1)}$  (so only the first block is zeroed out). Define  $\Pi_{(-1)}$  similarly. Then  $\sigma^\top B \sigma = \sigma_{(1)}^\top D_{z(1)} \Pi_{(1)}^\top \Pi_{(-1)}^\top D_{z(-1)} \sigma_{(-1)}$ . Define  $v := (\sigma_{(1)}^\top D_{z(1)} \Pi_{(1)}^\top \Pi_{(-1)}^\top D_{z(-1)})^\top$ . We first bound  $\|v\|_2$ . Note

$$\begin{aligned}
\|v\|_2 &= \sup_{\|y\|_2=1} v^\top y \\
&= \sigma_{(1)}^\top D_{z(1)} \Pi_{(1)}^\top \Pi_{(-1)}^\top D_{z(-1)} y \\
&= \sum_{j>1} \sigma_{(1)} D_{z(1)} \Pi_{(1)}^\top \Pi_{(j)}^\top D_{z(j)} y_{(j)} \\
&= \sum_{j>1} z_{(1)} D_{\sigma(1)} \Pi_{(1)}^\top \Pi_{(j)}^\top D_{z(j)} y_{(j)}
\end{aligned}$$



$$\leq \sum_{j>1} \|z_{(1)}\|_2 \cdot \|D_{\sigma(1)}\| \cdot \|\Pi_{(1)}^\top \Pi_{(j)}\| \|D_{z(j)}\| \|y_{(j)}\|_2$$

We have  $\|z_{(1)}\|_2 \leq \|z\|_2 = 1$  and  $\|D_\sigma\| = \|\sigma\|_\infty = 1$ . Also by  $2k$ -RIPness of  $\Pi$  and part (b) of [Lemma 5.3.6](#),  $\|\Pi_{(1)}^\top \Pi_{(j)}\| \leq c\varepsilon$  since  $j \neq 1$ . Also  $\|D_{z(j)}\| = \|z_{(j)}\|_\infty$ . Thus

$$\begin{aligned} v^\top y &\leq c\varepsilon \sum_{j>1} \|z_{(j)}\|_\infty \|y_{(j)}\|_2 \\ &\leq c\varepsilon \sum_{j>1} \left( \frac{\|z_{(j-1)}\|_1}{k} \right) \|y_{(j)}\|_2 && (5.10) \\ &\leq c\varepsilon \sum_{j>1} \left( \frac{\|z_{(j-1)}\|_2}{\sqrt{k}} \right) \|y_{(j)}\|_2 && (\text{Cauchy-Schwarz}) \\ &= \frac{c\varepsilon}{\sqrt{k}} \sum_{j>1} \|z_{(j-1)}\|_2 \cdot \|y_{(j)}\|_2 \\ &\leq \frac{c\varepsilon}{\sqrt{k}} \left( \sum_{j>1} \|z_{(j-1)}\|_2^2 \right)^{1/2} \left( \sum_{j>1} \|y_{(j)}\|_2^2 \right)^{1/2} && (\text{Cauchy-Schwarz}) \\ &\leq \frac{c\varepsilon}{\sqrt{k}} \|z\|_2 \cdot \|y\|_2 \\ &= \frac{c\varepsilon}{\sqrt{k}} \end{aligned}$$

The step [Eq. \(5.10\)](#) is known as “shelling”, i.e. dividing a vector into blocks after sorting entries by magnitude, then comparing some norm in the  $j$ th block with a norm in the  $(j-1)$ st block.

$$\begin{aligned} \mathbb{P}(|\sigma^\top B\sigma| > \varepsilon/4) &= \mathbb{P}(|v^\top \sigma| > \varepsilon/4) \\ &\leq 2e^{-C'\varepsilon^2/\|v\|_2^2} && (\text{Khinchine}) \\ &< 2^{-k} \end{aligned}$$

for  $c$  sufficiently small, which is at most  $\delta/4$ .

**Analyzing  $\mathcal{E}_3$ :**  $\sigma^\top B^\top \sigma = \sigma^\top B\sigma$ , so the analysis of this event is identical to that of  $\mathcal{E}_1$ .

**Analyzing  $\mathcal{E}_4$ :** We wish to bound  $\mathbb{P}_\sigma(|\sigma^\top C\sigma| > \frac{\varepsilon}{4})$ ; we will use the Hanson-Wright inequality ([Corollary 1.1.16](#)). For this we need to bound both  $\|C\|$  and  $\|C\|_F$ .

$$\begin{aligned} \|C\|_F^2 &= \sum_{\substack{i \neq j \\ i, j > 1}} \|D_{z(i)} \Pi_{(i)}^\top \Pi_{(j)} D_{z(j)}\|_F^2 \\ &\leq \sum_{i \neq j} \left( \|D_{z(i)}\| \cdot \|\Pi_{(i)}^\top \Pi_{(j)}\| \cdot \|D_{z(j)}\|_F \right)^2 && (\|AB\|_F \leq \|A\| \cdot \|B\|_F) \\ &\leq c^2 \varepsilon^2 \sum_{i \neq j} \|z(i)\|_\infty^2 \cdot \|z(j)\|_2^2 && (\text{Lemma 5.3.6}) \end{aligned}$$

$$\begin{aligned}
&\leq \frac{c^2 \varepsilon^2}{k} \sum_{i \neq j} \|z_{(i-1)}\|_2^2 \cdot \|z_{(j)}\|_2^2 && \text{(shelling)} \\
&\leq \frac{c^2 \varepsilon^2}{k} \left( \sum_i \|z_{(i)}\|_2^2 \right)^2 && \text{(monomials a superset of the terms in previous line)} \\
&= \frac{c^2 \varepsilon^2}{k}
\end{aligned}$$

To see that  $\|AB\|_F \leq \|A\| \|B\|_F$ , note that if the  $i$ th column of  $B$  is  $b^i$  then  $\|AB\|_F^2 = \sum_i \|Ab^i\|_2^2 \leq \|A\| \cdot \sum_i \|b^i\|_2^2 = \|A\|^2 \|B\|_F^2$ .

$$\begin{aligned}
\|C\| &= \sup_{\|y\|_2=1} |y^\top C y| \\
&= \left| \sum_{\substack{i \neq j \\ i, j > 1}} y_{(i)} D_{z_{(i)}} \Pi_{(i)}^\top \Pi_{(j)} D_{z_{(j)}} y_{(j)} \right| \\
&\leq \left| \sum_{\substack{i \neq j \\ i, j > 1}} \|y_{(i)}\|_2 \cdot \|D_{z_{(i)}}\| \cdot \|\Pi_{(i)}^\top \Pi_{(j)}\| \cdot \|D_{z_{(j)}}\| \cdot \|y_{(j)}\|_2 \right| \\
&\leq \frac{c\varepsilon}{k} \sum_{\substack{i \neq j \\ i, j > 1}} \|y_{(i)}\|_2 \cdot \|z_{(i-1)}\|_2 \cdot \|z_{(j-1)}\|_2 \cdot \|y_{(j)}\|_2 && \text{(shelling twice and Lemma 5.3.6)} \\
&\leq \frac{c\varepsilon}{k} \left( \sum_{i>1} \|y_{(i)}\|_2 \cdot \|z_{(i-1)}\|_2 \right)^2 \\
&\leq \frac{c\varepsilon}{k} \left( \sum_{i>1} \|y_{(i)}\|_2^2 \right) \cdot \left( \sum_{i>1} \|z_{(i-1)}\|_2^2 \right) && \text{(Cauchy-Schwarz)} \\
&\leq \frac{c\varepsilon}{k} \|y\|_2^2 \cdot \|z\|_2^2 \\
&= \frac{c\varepsilon}{k}.
\end{aligned}$$

We thus have by Cauchy-Schwarz that

$$\mathbb{P}(|\sigma^\top C \sigma| > \frac{\varepsilon}{4}) \lesssim e^{-C \frac{\varepsilon^2}{16} \cdot \frac{k}{c^2 \varepsilon^2}} + e^{-C \frac{\varepsilon}{4} \cdot \frac{k}{c\varepsilon}},$$

which is at most  $\delta/4$  for  $c$  sufficiently small. □

# Chapter 6

## Linear algebra applications

In this chapter we focus on applications of sketching to linear algebra problems such as approximate matrix multiplication, regression, low-rank approximation, and  $k$ -means clustering (which is actually a special case of constrained low-rank approximation). A more in-depth treatment of the use of sketching in linear algebra applications can be found in the book [Woo14].

### 6.1 Approximate matrix multiplication

Suppose we have two matrices  $A \in \mathbb{R}^{n \times p}, B \in \mathbb{R}^{n \times d}$  with  $n$  large. We could naively compute  $A^\top B$  using for loops, which would take  $O(ndp)$  flops. Asymptotically faster matrix multiplication algorithms do exist (the current records are in [AW20, GU18]), but are of only theoretical interest as they are slow in practice. We denote the rows of  $A, B$  as  $a_i, b_i$ , respectively:

$$A = \begin{pmatrix} \text{---} & a_1^\top & \text{---} \\ \text{---} & a_2^\top & \text{---} \\ & \vdots & \\ & \vdots & \\ & \vdots & \\ & \vdots & \\ & \vdots & \\ \text{---} & a_n^\top & \text{---} \end{pmatrix}, \quad B = \begin{pmatrix} \text{---} & b_1^\top & \text{---} \\ \text{---} & b_2^\top & \text{---} \\ & \vdots & \\ & \vdots & \\ & \vdots & \\ & \vdots & \\ & \vdots & \\ \text{---} & b_n^\top & \text{---} \end{pmatrix}$$

There are two main approaches to using sketching to speed up the computation of some matrix  $C$  that approximates  $A^\top B$ : a sampling approach [DKM06], and an oblivious linear sketching approach [Sar06] (where “oblivious” refers to the fact that the sketching matrix  $\Pi$  is chosen without knowledge of the inputs  $A, B$ ). Specifically, for some norm  $\|\cdot\|_X$ , we would like that  $\|A^\top B - C\|_X$  is small with good probability over the randomness used by the algorithm. We focus here on  $\|\cdot\|_X = \|\cdot\|_F$  being the Frobenius norm.

#### 6.1.1 Sampling approach

Here we describe the approach of [DKM06] for approximate matrix multiplication based on row-sampling. The main starting point of the approach is the identity

$$A^\top B = \sum_{i=1}^n a_i b_i^\top.$$

That is, the matrix we want  $A^\top B$  can be represented as the sum of  $n$  rank-1 matrices. Computing one such rank-1 matrix and adding it into a running sum takes  $O(dp)$  time, so computing all  $n$  of them takes  $O(ndp)$  time. The idea is then to approximate  $A^\top B$  by only sampling  $m$  of these matrices, thus improving the runtime to only  $O(mdp)$ . That is, we pick some probability distribution  $\vec{p} = (p_1, \dots, p_n)$  on the rows then pick  $m$  i.i.d. rows  $i_1, \dots, i_m$  from distribution  $\vec{p}$ . We then approximate  $A^\top B$  as

$$C := \frac{1}{m} \sum_{r=1}^m \frac{a_{i_r} b_{i_r}^\top}{p_{i_r}}.$$

We note that  $C$  can also be seen as  $(\Pi A)^\top (\Pi B)$  for a linear sketching matrix  $\Pi \in \mathbb{R}^{m \times n}$  which is *not* oblivious, i.e. the distribution from which we draw  $\Pi$  depends on  $A, B$ . Specifically,  $\Pi$  is a sampling matrix where the  $r$ th row has a  $1/\sqrt{m}$  in the  $i_r$ th position and zeroes elsewhere.

One can show that  $\mathbb{E} C = A^\top B$  regardless of how the  $p_i$  are chosen (as long as none of them is zero). The question then is how to pick  $\vec{p}$  so that not only is  $C$  an unbiased estimator of  $A^\top B$ , but also  $C$  is close to  $A^\top B$  with good probability. Specifically, by applying Markov after squaring both sides of the inequality,

$$\mathbb{P}_{\Pi}(\|(\Pi A)^\top (\Pi B) - A^\top B\|_F > \varepsilon \|A\|_F \|B\|_F) < \frac{\mathbb{E} \|(\Pi A)^\top (\Pi B) - A^\top B\|_F^2}{\varepsilon^2 \|A\|_F^2 \|B\|_F^2}. \quad (6.1)$$

**Lemma 6.1.1.** *Set  $p_i = \frac{\|a_i\|_2 \cdot \|b_i\|_2}{\sum_{j=1}^n \|a_j\|_2 \cdot \|b_j\|_2}$  and  $m \geq 1/(\varepsilon^2 \eta)$ . Then for  $\varepsilon, \eta \in (0, 1/2)$ ,*

$$\mathbb{P}_{\Pi}(\|(\Pi A)^\top (\Pi B) - A^\top B\|_F > \varepsilon \|A\|_F \|B\|_F) < \eta.$$

*Proof.* By Eq. (6.1), it suffices to show  $\mathbb{E} \|(\Pi A)^\top (\Pi B) - A^\top B\|_F^2 \leq \|A\|_F^2 \|B\|_F^2 / m$ . Define  $Z_r = \frac{1}{m} \frac{a_{i_r} b_{i_r}^\top}{p_{i_r}}$  so that  $(\Pi A)^\top (\Pi B) = \sum_{r=1}^m Z_r$ . Then

$$\begin{aligned} \mathbb{E} \|(\Pi A)^\top (\Pi B) - A^\top B\|_F^2 &= \sum_{i,j} \mathbb{E} \left( \sum_{r=1}^m (Z_r - \mathbb{E} Z_r)_{i,j} \right)^2 \\ &= \sum_{i,j} \text{Var} \left[ \sum_{r=1}^m (Z_r)_{i,j} \right] \\ &= \sum_{i,j} \sum_{r=1}^m \text{Var}[(Z_r)_{i,j}] \quad (\text{independence of the } Z_r) \end{aligned}$$

Now, the  $Z_r$  are distributed as a random variable  $Z$  where

$$Z_{i,j} = \frac{1}{m} \sum_{k=1}^n \rho_k \frac{(a_k)_i (b_k)_j}{p_k}$$

where  $\rho_k$  is an indicator random variable for the event that row  $k$  was sampled. Thus

$$\begin{aligned} \text{Var}[Z_{i,j}] &\leq \mathbb{E} Z_{i,j}^2 \\ &= \frac{1}{m^2} \sum_{k=1}^n p_k \frac{(a_k)_i^2 (b_k)_j^2}{p_k^2} \\ &= \frac{1}{m^2} \sum_{k=1}^n \frac{(a_k)_i^2 (b_k)_j^2}{\|a_k\|_2 \|b_k\|_2} \cdot \left( \sum_{t=1}^n \|a_t\|_2 \|b_t\|_2 \right). \end{aligned}$$

Plugging back into Eq. (independence of the  $Z_r$ ),

$$\begin{aligned}
\mathbb{E} \|(\Pi A)^\top (\Pi B) - A^\top B\|_F^2 &\leq \frac{1}{m} \left( \sum_{t=1}^n \|a_t\|_2 \|b_t\|_2 \right) \cdot \sum_k \left( \sum_{i,j} \frac{(a_k)_i^2 (b_k)_j^2}{\|a_k\|_2 \|b_k\|_2} \right) \\
&= \frac{1}{m} \left( \sum_{t=1}^n \|a_t\|_2 \|b_t\|_2 \right) \cdot \left( \sum_k \|a_k\|_2 \|b_k\|_2 \right) \\
&= \frac{1}{m} \left( \sum_{t=1}^n \|a_t\|_2 \|b_t\|_2 \right)^2 \\
&\leq \frac{1}{m} \left( \sum_{t=1}^n \|a_t\|_2^2 \right) \left( \sum_{t=1}^n \|b_t\|_2^2 \right) \quad (\text{Cauchy-Schwarz}) \\
&= \frac{\|A\|_F^2 \|B\|_F^2}{m}
\end{aligned}$$

as desired.  $\square$

One downside of Lemma 6.1.1 is that the dependence on the failure probability is  $1/\eta$ . For failure probability  $\delta \ll 1$ , we would prefer to have  $O(\log(1/\delta))$  dependence. Our usual method for accomplishing would be to instantiate the above scheme for  $\eta = 1/4$ , then run the above scheme  $R = \Theta(\log(1/\delta))$  times to get approximations  $C_1, \dots, C_r$ , then return the “median” of the  $C_r$  as our output. The issue here is that unlike previous problems where the output was a number, here the outputs are matrices so it is not clear what the “median” of a set of matrix even means. One way around this issue is the following, suggested by [CW09].

Let  $C^* = A^\top B$  be the true value of the matrix product. The Chernoff bound tells us that with probability at least  $1 - \exp(-\Omega(R))$ , strictly more than  $2R/3$  of the  $C_r$  satisfy  $\|C_r - C^*\|_F \leq \gamma$ , where  $\gamma = \varepsilon \|A\|_F \|B\|_F$  is our desired error bound. Let  $S \subset [R]$  be this “good” subset of the indices  $r$ . Then  $|S| > 2R/3$ . We also know that if  $r, r' \in R$ , then by the triangle inequality  $\|C_r - C_{r'}\|_F \leq \|C_r - C^*\|_F + \|C_{r'} - C^*\|_F \leq 2\gamma$ . Meanwhile if  $C_{r'}$  is “far” from a good approximation, i.e.  $\|C_{r'} - C^*\|_F > 3\gamma$ , then by triangle inequality we have that for any  $r \in S$ ,  $\|C_{r'} - C_r\|_F \geq \|C_{r'} - C^*\|_F - \|C_r - C^*\|_F > 2\gamma$ . This motivates then the scheme of [CW09]: loop over all  $r \in [R]$  and return the first  $C_r$  for which  $|\{r' \in [R] : \|C_r - C_{r'}\|_F \leq 2\gamma\}| \geq 2R/3$ . By the given analysis, any such  $r$  is guaranteed to have  $\|C_r - C^*\|_F \leq 3\gamma$ , and such an  $r$  must exist since any  $r \in S$  would do.

The downside of the above scheme is that it requires  $\Theta(R^2)$  distance comparisons, whereas computing the median in 1 dimension only required  $O(R)$  comparisons. One way around this is to use a randomized Las Vegas scheme: instead of trying all  $r$  in some fixed order, pick a random  $r$  and test whether it satisfies  $|\{r' \in [R] : \|C_r - C_{r'}\|_F \leq 2\gamma\}| \geq 2R/3$ . Since there are at least  $2R/3$  values of  $r$  which do (namely any  $r \in S$ ), the probability of a good  $R$  is at least  $2/3$ , and thus the expected number of  $r$  we have to test is thus only  $O(1)$ . Therefore the expected number of distance comparisons is only  $O(R)$ . Alternatively, Narayanan showed a *deterministic* algorithm which is always guaranteed to compute only a linear number of distances to accomplish this task [Nar18].

### 6.1.2 Oblivious linear sketching approach

An *oblivious* linear sketching approach is one where our linear sketch matrix  $\Pi$  does not depend on  $A, B$ . In particular it will be chosen randomly from some distribution  $\mathcal{D}$  which does not depend on  $A, B$ . We will use such a distribution that satisfies what is known as the *JL moment property*.

**Definition 6.1.2** ([KN14]). We say a distribution  $\mathcal{D}$  over  $\mathbb{R}^{m \times n}$  satisfies the  $(\varepsilon, \delta, p)$ -JL moment property if for all unit vectors  $z \in \mathbb{R}^n$ ,

$$\| \|\Pi z\|_2^2 - 1 \|_p \leq \varepsilon \delta^{1/p}.$$

The above definition is useful for obtaining DJL since, by Markov's inequality after raising both sides to the  $p$ th power,

$$\mathbb{P}_{\Pi \sim \mathcal{D}} (|\|\Pi z\|_2^2 - 1| > \varepsilon) < \frac{\| \|\Pi z\|_2^2 - 1 \|_p^p}{\varepsilon^p}.$$

The above is at most  $\delta$  if  $\mathcal{D}$  satisfies the  $(\varepsilon, \delta, p)$ -JL moment property.

Though we did not state it explicitly in previous chapters, several distributions  $\mathcal{D}$  we have seen satisfy the JL moment property. For example,  $\Pi$  could be a CountSketch with  $m = O(1/(\varepsilon^2 \delta))$ , which satisfies the property for  $p = 2$ . Alternatively a matrix of i.i.d. random signs, or Sparse JL transform, with  $m = O(\varepsilon^{-2} \log(1/\delta))$  (and  $s = O(\varepsilon^{-1} \log(1/\delta))$  for SparseJL) satisfy the JL moment property for  $p = \Theta(\log(1/\delta))$  (this follows by inspecting which moment is being used in the Hanson-Wright inequality in those analyses). It can also be shown that the FastJL transform of [AC09] satisfies the JL moment property (see [CNW16]).

We now show that any distribution with the JL moment property for  $p \geq 2$  provides an approximate matrix multiplication guarantee. First we give a lemma, which is based on one in [KN14] but with an argument that gives an improved constant factor as observed by Olivier Zahm.

**Lemma 6.1.3.** *Suppose  $x, y \in \mathbb{R}^n$  have  $\|x\|_2 = \|y\|_2 = 1$  and suppose  $\mathcal{D}$  satisfies the  $(\varepsilon, \delta, p)$ -JL moment property. Then for  $\Pi \sim \mathcal{D}$ ,*

$$\| \langle \Pi x, \Pi y \rangle - \langle x, y \rangle \|_p \leq \varepsilon \delta^{1/p}.$$

*Proof.* Using the relation  $\langle x, y \rangle = \frac{1}{4}(\|x + y\|_2^2 - \|x - y\|_2^2)$ ,

$$\begin{aligned} \| \langle \Pi x, \Pi y \rangle - \langle x, y \rangle \|_p &= \frac{1}{4} \| (\|\Pi(x + y)\|_2^2 - \|x + y\|_2^2) - (\|\Pi(x - y)\|_2^2 - \|x - y\|_2^2) \|_p \\ &\leq \frac{1}{4} [ \| \|\Pi(x + y)\|_2^2 - \|x + y\|_2^2 \|_p + \| \|\Pi(x - y)\|_2^2 - \|x - y\|_2^2 \|_p ] \\ &\hspace{15em} (\text{triangle inequality}) \\ &\leq \frac{\|x + y\|_2^2}{4} \cdot \| \|\Pi(\frac{x + y}{\|x + y\|_2})\|_2^2 - \|\frac{x + y}{\|x + y\|_2}\|_2^2 \|_p \\ &\quad + \frac{\|x - y\|_2^2}{4} \cdot \| \|\Pi(\frac{x - y}{\|x - y\|_2})\|_2^2 - \|\frac{x - y}{\|x - y\|_2}\|_2^2 \|_p \\ &\leq \frac{\varepsilon \delta^{1/p}}{4} \cdot (\|x + y\|_2^2 + \|x - y\|_2^2) \hspace{5em} (\text{JL moment property}) \\ &= \frac{\varepsilon \delta^{1/p}}{4} \cdot (2\|x\|_2^2 + 2\|y\|_2^2) \\ &= \varepsilon \delta^{1/p}. \end{aligned}$$

□

**Theorem 6.1.4.** *Suppose  $\Pi \in \mathbb{R}^{m \times n}$  is drawn from a distribution satisfying the  $(\varepsilon, \delta, p)$ -JL moment property for some  $p \geq 2$ . Then*

$$\mathbb{P}_{\Pi} \left( \| (\Pi A)^\top (\Pi B) - A^\top B \|_F > \varepsilon \|A\|_F \|B\|_F \right) < \delta.$$

*Proof.* We will use Markov's inequality, which gives

$$\mathbb{P}_{\Pi} \left( \|(\Pi A)^{\top}(\Pi B) - A^{\top}B\|_F > \varepsilon \|A\|_F \|B\|_F \right) < \frac{\|(\Pi A)^{\top}(\Pi B) - A^{\top}B\|_F^p}{\varepsilon^p \|A\|_F^p \|B\|_F^p}$$

Let  $a^i$  denote the  $i$ th column of  $A$ , and similarly  $b^i$  is the  $i$ th column of  $B$ . Define  $\tilde{a}^i = a^i / \|a^i\|_2$ , and similarly for  $\tilde{b}^i$ . Write  $M := (\Pi A)^{\top}(\Pi B) - A^{\top}B$ . Then

$$M_{i,j} = (\langle \Pi \tilde{a}^i, \Pi \tilde{b}^j \rangle - \langle \tilde{a}^i, \tilde{b}^j \rangle) \|a^i\|_2 \|b^j\|_2.$$

It thus suffices to show  $\|M\|_F \leq \varepsilon \|A\|_F \|B\|_F \delta^{1/p}$ . We have

$$\begin{aligned} \|M\|_F^p &= \|M\|_F^2 \|M\|_F^{p-2} \\ &= \left\| \sum_{i,j} \underbrace{(\langle \Pi \tilde{a}^i, \Pi \tilde{b}^j \rangle - \langle \tilde{a}^i, \tilde{b}^j \rangle)^2}_{X_{i,j}} \|a^i\|_2^2 \|b^j\|_2^2 \right\|_F^{p/2} \\ &\leq \left( \sum_{i,j} \|a^i\|_2^2 \cdot \|b^j\|_2^2 \cdot \|X_{i,j}\|_{p/2}^2 \right)^{1/2} && \text{(triangle inequality, since } p/2 \geq 1) \\ &\leq \sum_{i,j} \|a^i\|_2 \cdot \|b^j\|_2 \cdot \|X_{i,j}\|_{p/2}^{1/2} && (\sqrt{A+B} \leq \sqrt{A} + \sqrt{B}) \\ &\leq \sum_{i,j} \|a^i\|_2 \cdot \|b^j\|_2 \cdot \|X_{i,j}\|_p \\ &\leq \varepsilon \delta^{1/p} \left( \sum_{i,j} \|a^i\|_2 \cdot \|b^j\|_2 \right) && \text{(Lemma 6.1.3)} \\ &\leq \varepsilon \delta^{1/p} \left( \sum_i \|a^i\|_2^2 \right)^{1/2} \cdot \left( \sum_i \|b^i\|_2^2 \right)^{1/2} && \text{(Cauchy-Schwarz)} \\ &= \varepsilon \delta^{1/p} \|A\|_F \|B\|_F. \end{aligned}$$

□

As mentioned above, [Theorem 6.1.4](#) implies we can get approximate matrix multiplication guarantees using the CountSketch matrix, SparseJL matrices, dense Rademacher-entried JL matrices, or FastJL transforms.

## 6.2 Subspace embeddings

The notion of a *subspace embedding* was introduced by Sarlós [[Sar06](#)] and is related (but not identical to) approximate matrix multiplication with respect to the  $\ell_2 \rightarrow \ell_2$  operator norm  $\|\cdot\|$ . Specifically, suppose we wanted an AMM guarantee with respect to this norm for multiplying  $A^{\top}A$ . Then we would want

$$\|(\Pi A)^{\top}(\Pi A) - A^{\top}A\| < \varepsilon \|A\|^2 = \varepsilon \|A^{\top}A\|. \quad (6.2)$$

Using that for real symmetric matrices  $M$  we have  $\|M\| = \sup_{\|x\|_2=1} |x^{\top}Mx|$ , Thus the guarantee [Eq. \(6.2\)](#) would be equivalent to

$$\forall \|x\|_2 = 1, \quad \left| \|\Pi Ax\|_2^2 - \|Ax\|_2^2 \right| \leq \sup_{\|z\|_2=1} \varepsilon \|Az\|_2^2.$$

A subspace embedding (for the column space of  $A$ ) is instead  $\Pi$  satisfying the stronger guarantee

$$\forall \|x\|_2 = 1, \quad \left| \|\Pi Ax\|_2^2 - \|Ax\|_2^2 \right| \leq \varepsilon \|Ax\|_2^2.$$

**Definition 6.2.1.** Let  $E \subset \mathbb{R}^n$  be a linear subspace. Then  $\Pi \in \mathbb{R}^{m \times n}$  is an  $\varepsilon$ -subspace embedding for  $E$  if

$$\forall x \in E, \quad (1 - \varepsilon)\|x\|_2^2 \leq \|\Pi x\|_2^2 \leq (1 + \varepsilon)\|x\|_2^2.$$

Let  $U \in \mathbb{R}^{n \times d}$  be a matrix with orthonormal columns that form a basis for  $E$ , so that  $E = \{Uz : z \in \mathbb{R}^d\}$ . Then being a subspace embedding is equivalent to the condition

$$\forall \|z\|_2 = 1, \quad \left| \|\Pi Uz\|_2^2 - \|Uz\|_2^2 \right| \leq \varepsilon \|Uz\|_2^2.$$

Note  $\|Uz\|_2^2 = z^\top U^\top Uz = \|z\|_2^2$ , and thus the above is equivalent to

$$\sup_{\|z\|_2=1} \left| \|\Pi Uz\|_2^2 - 1 \right| \leq \varepsilon,$$

i.e.  $\|(\Pi U)^\top (\Pi U) - I\| \leq \varepsilon$ . We will often make use of this alternate, equivalent definition of subspace embeddings. Note this definition is a natural generalization of the condition in the DJL guarantee, where  $|\cdot|$  is replaced by  $\|\cdot\|$ , and the unit vector  $u$  in DJL (i.e. a vector  $u \in \mathbb{R}^{n \times 1}$  with  $u^\top u = 1$ ) is being replaced by a matrix  $U$  satisfying a similar, higher-dimensional version  $U^\top U = I$ .

It is known that subspace embeddings can be used to design algorithms for approximate computational problems in linear algebra, such as approximating leverage scores, least squares regression, low-rank approximation,  $k$ -means clustering, and several others. We will go into a subset of these in future sections. For the remainder of this section, we discuss how one obtains subspace embeddings.

### 6.2.1 Given an orthonormal basis

If we know (or compute) an orthonormal basis  $U \in \mathbb{R}^{n \times d}$  for the subspace  $E$  of interest (i.e.  $U^\top U = I$ ), then we have that  $\Pi = U^\top \in \mathbb{R}^{d \times d}$  is an  $\varepsilon$ -subspace embedding for  $A$  with  $\varepsilon = 0$ . This is because any  $x \in E$  can be written as  $Uz$  so that

$$\|x\|_2^2 = \|Uz\|_2^2 = \|z\|_2^2$$

But then  $\|U^\top x\|_2^2 = \|U^\top Uz\|_2^2 = \|z\|_2^2$  as well. Furthermore note that a subspace embedding can never have fewer than  $d$  rows, since otherwise some non-zero vector in  $E$  will be in the kernel of  $\Pi$  and thus not have its norm preserved at all.

### 6.2.2 Leverage score sampling

A sampling approach for obtaining a subspace embedding was first proposed in [DMM06] then further developed and analyzed in [SS11]. The idea is to again use the fact that matrix-matrix product can be represented as the sum of rank-1 outerproducts then sample, as in Subsection 6.1.1:

$$A^\top A = \sum_{i=1}^n a_i a_i^\top.$$

Suppose we keep each row  $i$  independently with probability  $p_i$ . That is, we set  $A = \sum_{i=1}^n \eta_i \frac{e_i a_i^\top}{\sqrt{p_i}}$  where  $\eta_i$  be an indicator random variable for the event that we do not zero out the  $i$ th row (so the  $i$ th row of  $\Pi A$  is  $a_i^\top / \sqrt{p_i}$  with probability  $p_i$ , and is 0 otherwise). Then

$$(\Pi A)^\top (\Pi A) = \sum_{i=1}^n \eta_i \cdot \frac{a_i a_i^\top}{p_i}.$$



In the case of AMM with respect to Frobenius norm, we set the  $p_i$  proportional to  $\|a_i\|_2 \cdot \|b_i\|_2$ , but what should they be set to in order to guarantee the subspace embedding property with good probability? Specifically, we would like

$$\mathbb{P}_{\Pi} \left( \|(\Pi U)^\top (\Pi U) - I\| > \varepsilon \right) < \delta$$

while keeping  $\sum_{i=1}^n p_i$  small (since that is the expected target dimension).

To guide us in selecting the  $p_i$ , we make the following definition:

**Definition 6.2.2.** Given a matrix  $A \in \mathbb{R}^{n \times d}$  with rows  $a_i^\top$ , we define the  $i$ th *sensitivity*  $R_i$  by

$$R_i := \sup_{\|x\|_2=1} \frac{\langle a_i, x \rangle^2}{\|Ax\|_2^2}.$$

Note the denominator equals  $\sum_j \langle a_j, x \rangle^2$ , and thus  $0 \leq R_i \leq 1$  always.

Now, let us build some intuition: it would be strange to set any  $p_i = 0$  (since then we are *always* ignoring some part of the input matrix  $A$ ). Then if  $p_i > 0$ , for the  $x$  achieving the sup in the definition of sensitivity, any time we do sample row  $i$  we would be guaranteed to have  $\|\Pi Ax\|_2^2 \geq (1/p_i) \cdot \langle a_i, x \rangle^2 = (R_i/p_i) \|Ax\|_2^2$ . Thus if we want to preserve  $\|Ax\|_2^2$  multiplicatively, we must certainly have  $p_i \geq R_i/2$  (else we alter the squared norm by a factor of more than 2). We now give an alternative definition. First, we introduce some more linear algebra background:

**Definition 6.2.3.** For a matrix  $A \in \mathbb{R}^{n \times d}$  of rank  $r \leq d$ , the *singular value decomposition (SVD)* of  $A$  is

$$A = U \Sigma V^\top$$

where  $U, V$  each have  $r$  orthonormal columns, and  $\Sigma$  is a diagonal matrix with  $\sigma_i := \Sigma_{i,i}$  with  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ .

It is a theorem in linear algebra, which we will not prove here, that every real matrix has a (not necessarily unique) SVD (though  $\Sigma$  is unique). We then have that the orthogonal projection onto the column space of  $X$  is  $UU^\top$ . We also have the following definition.

**Definition 6.2.4.** For a matrix with SVD decomposition  $A = U \Sigma V^\top$ , the *Moore-Penrose pseudoinverse* is defined to be  $A^+ := V \Sigma^{-1} U^\top$ . Then  $AA^+ = UU^\top$  is the orthogonal projection onto the column space of  $A$ , and  $A^+A = VV^\top$  is the orthogonal projection onto the rowspace of  $A$ .

**Definition 6.2.5.** Given a matrix  $A \in \mathbb{R}^{n \times d}$  with rows  $a_i^\top$ , we define the  $i$ th *leverage score*  $\ell_i$  by

$$\ell_i := a_i^\top (A^\top A)^+ a_i = e_i^\top A (A^\top A)^+ A^\top e_i.$$

Since  $UU^\top$  is the orthogonal projection onto the column space of  $A$ , equivalently the  $i$ th leverage score is the squared norm of the  $i$ th row of  $U$ .

We now show that the definitions ‘‘sensitivity’’ and ‘‘leverage score’’ are in fact equivalent.

**Lemma 6.2.6.** *Let  $R_i$  be the sensitivity of the  $i$ th row and  $\ell_i$  be its leverage score. Then  $R_i = \ell_i$ .*

*Proof.* Note that if  $M$  is invertible,  $AM$  and  $A$  have the same sensitivities (since the new maximizer is simply  $M^{-1}x$ ) and leverage scores (since the column spaces of  $AM$  and  $A$  are identical). Choose  $M$  to be a matrix such that  $AM$  has orthonormal columns, i.e.  $M = V \Sigma^{-1}$  (taken from the SVD). Thus  $AM = U$ , with rows  $u_1, \dots, u_n$ . Then  $\ell_i = \|u_i\|_2^2$ , whereas

$$R_i = \sup_{\|x\|_2=1} \frac{\langle u_i, x \rangle^2}{\|Ux\|_2^2}$$

reference	sketch name	# rows $m$	mult. time
folklore	gaussian matrix	$\varepsilon^{-2}(d + \log(1/\delta))$	$mnd$
[CW17, NN13a, MM13]	CountSketch	$\varepsilon^{-2}\delta^{-1}d^2$	$nnz(A) + md$
[NN13a, BDN15, Coh16]	OSNAP	$\varepsilon^{-2}d \log(d/\delta)$	$\varepsilon^{-1} nnz(A) \log(d/\delta) + md$
[Sar06, Tro11, CNW16]	SRHT	$\varepsilon^{-2} \log(d/\delta)(d + \log(1/(\varepsilon\delta)))$	$nd \log n$
[SS11]	leverage score sampling	$\varepsilon^{-2}d \log(d/\delta)$	$md$
[BSS12, LS17, LS18]	row subset selection	$d/\varepsilon^2$	$md$

Figure 6.1: The last column is the time it take to compute  $\Pi A$  given  $A \in \mathbb{R}^{n \times d}$ , where we want  $\Pi \in \mathbb{R}^{m \times n}$  to be an  $\varepsilon$ -subspace embedding for the column space of  $A$ . The first four rows are OSE's succeeding with probability  $1 - \delta$  ("SRHT" stands for "Subsampled Randomized Hadamard Transform"), where  $\Pi = \frac{1}{\sqrt{m}} SHD$  as in Eq. (5.9). The last two are not oblivious, but rather select a subset of  $m$  rows (possibly rescaled) of  $A$ . Leverage score sampling is discussed in Subsection 6.2.2. Row subset selection (also called BSS or column subset selection) is an algorithm that selects the  $m$  rows given  $A$ ; the more interesting point for these constructions is not the time to multiply  $\Pi A$  (last column), but rather the time it takes to select the rows.

$$= \sup_{\|x\|_2=1} \langle u_i, x \rangle^2. \quad (\|Ux\|_2^2 = \|x\|_2^2 = 1)$$

The supremum is achieved for  $x = u_i / \|u_i\|_2$ , which gives  $R_i = \|u_i\|_2^2 = \ell_i$ .  $\square$

Thus we should pick, as a lower bound, at least  $p_i \geq \ell_i/2$ . Note than  $\sum_i p_i = \Omega(\sum_i \ell_i)$ . We have  $\sum_i \ell_i = \|U\|_F^2$ , which is exactly the rank of  $A$  (since the squared Frobenius norm is the sum of the squared column norms of  $U$ , which are each exactly 1). A theorem of Spielman and Srivastava, based on the *Matrix Chernoff bound* (one can also use the *non-commutative Khintchine inequality*) shows that in fact one can actually take  $p_i \approx \ell_i$  and obtain good results. We do not show the proof here.

**Theorem 6.2.7** ([SS11]). *Suppose  $A \in \mathbb{R}^{n \times d}$  has rank at most  $r$ . For some universal constant  $C > 0$  and any  $\varepsilon, \delta \in (0, 1/2)$ , choose  $p_i := \min\{1, C\varepsilon^{-2}\ell_i \log(r/\delta)\}$ . Then for  $\Pi$  as described above,*

$$\mathbb{P}_{\Pi}(\|(\Pi U)^{\top}(\Pi U) - I\| > \varepsilon) < \delta.$$

### 6.2.3 Oblivious subspace embeddings

The downside of leverage score sampling is that given some arbitrary matrix  $A$ , it is not clear what its leverage scores are. They can be computed via the SVD, but that is slow. Alternative fast approximation algorithms for the leverage scores do exist though [DMMW12].

Alternatively, one could use a distribution  $\mathcal{D}$  over  $\Pi$  which does not depend on the input matrix  $A$ . Such a distribution is called an *oblivious subspace embedding*, introduced by Sarlós [Sar06].

**Definition 6.2.8.** An  $(\varepsilon, d, \delta)$ -oblivious subspace embedding (OSE), is a distribution  $\mathcal{D}$  over  $\mathbb{R}^{m \times n}$  such that for any matrix  $U \in \mathbb{R}^{n \times d}$  with orthonormal columns (i.e. for any linear subspace, which is the column space of  $U$ ),

$$\mathbb{P}_{\Pi \sim \mathcal{D}}(\|(\Pi U)^{\top}(\Pi U) - I\| > \varepsilon) < \delta.$$

There are generally three ways to show that  $\mathcal{D}$  a distribution over  $\mathbb{R}^{m \times n}$  is an OSE.

**Net argument.** For a subspace  $E$  of dimension  $d$ , let  $E'$  be a  $\gamma$ -net of unit Euclidean ball in the subspace for  $\gamma$  a sufficiently small constant (like  $1/4$ ). Thus  $|E'| \leq O(1/\gamma)^d$  by Lemma 5.2.4. Then as long as  $\Pi$  preserves the norms of all  $x' \in E'$  up to  $1 \pm \varepsilon$ , it is possible to show that  $\Pi$  then preserves *all* vectors in  $E$  up to  $1 \pm O(\varepsilon)$  (see [CW17]). Thus we can e.g. have  $M = O(\log |E'|/\varepsilon^2) = O(d/\varepsilon^2)$  by the JL lemma, where the entries in  $\Pi$  have i.i.d.  $\pm 1/\sqrt{m}$  entries. It is known that such  $m$  is optimal for OSE's [NN14].

**Moment method.** By Markov's inequality, for any  $p$  we have

$$\mathbb{P}_{\Pi}(\|(\Pi U)^{\top}(\Pi U) - I\| > \varepsilon) < \frac{\mathbb{E} \|(\Pi U)^{\top}(\Pi U) - I\|^p}{\varepsilon^p}$$

Write  $M := (\Pi U)^{\top}(\Pi U) - I$ . Then  $M$  is a real symmetric matrix and thus by the spectral theorem has all real eigenvalues  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_d|$ . Thus  $\|M\| = |\lambda_1|$ , which equals  $\lambda_1^p$  if  $p$  is an even integer. Meanwhile,  $\text{tr}(M^p) = \sum_i \lambda_i^p$ . Thus for  $p$  an even integer,  $\|M\|^p \leq \text{tr}(M^p)$ , and thus

$$\mathbb{P}_{\Pi}(\|(\Pi U)^{\top}(\Pi U) - I\| > \varepsilon) < \frac{\mathbb{E} \text{tr}(((\Pi U)^{\top}(\Pi U) - I)^p)}{\varepsilon^p}.$$

Next, it can be shown by induction on  $p$  that for any square matrix  $M$  and integer  $p > 0$ ,

$$(M^p)_{i,j} = \sum_{i_1=i, i_2, \dots, i_{p+1}=j} \prod_{t=1}^p M_{i_t, i_{t+1}}.$$

Therefore

$$\text{tr}(M^p) = \sum_{\substack{i_1, \dots, i_{p+1} \\ i_1=i_{p+1}}} \prod_{t=1}^p M_{i_t, i_{t+1}},$$

whose expectation we can then compute using linearity of expectation (see for example [NN13a]). Alternatively, for some OSE's it is possible to provide simpler proofs using matrix exponential based arguments, similar to the MGF-based approach to prove the Chernoff bound (see for example [Coh16]).

**Approximate matrix multiplication.** Since  $\|M\|_F \leq \|M\|$  (the LHS is the  $\ell_2$  norm of the singular values, and the RHS is their  $\ell_\infty$  norm), we have

$$\|(\Pi U)^{\top}(\Pi U) - I\| \leq \|(\Pi U)^{\top}(\Pi U) - I\|_F,$$

thus as long as the RHS is at most  $\varepsilon$ ,  $\Pi$  is an  $\varepsilon$ -subspace embedding. The approximate matrix multiplication guarantee of Subsection 6.1.2 bounds the above by  $\varepsilon \|U\|_F^2 = \varepsilon d$ . We can thus apply AMM with error parameter  $\varepsilon' = \varepsilon/d$  to obtain a subspace embedding. This for example shows that the CountSketch yields an OSE with  $m = O(d^2/(\varepsilon^2\delta))$  (see also [CW17, NN13a, MM13]). Note this result is simply syntactic sugar for the moment method of the previous section, since ultimately the AMM property for the CountSketch follows from bounding  $\mathbb{E} \|M\|_F^2$  for  $M = (\Pi U)^{\top}(\Pi U) - I$ . This is the same as  $\mathbb{E} \text{tr}(M^2)$  since  $\|M\|_F^2$  is generally equal to  $\text{tr}(M^{\top}M)$  for any real matrix  $M$ , which equals  $\text{tr}(M^2)$  if  $M$  is symmetric. Having such sparse  $M$  is advantageous since then for any input matrix  $A$ , we can compute  $\Pi A$  in time proportional to  $m$  plus the number of nonzero entries in  $A$ .

### 6.3 Least squares regression

In least squares regression we are given as input  $X \in \mathbb{R}^{n \times d}$ ,  $y \in \mathbb{R}^n$  and would like to compute

$$\beta^{LS} = \underset{\beta}{\operatorname{argmin}} \|X\beta - y\|_2.$$

Note  $y$  can be decomposed into  $y^\perp + y^\parallel$ , where  $y^\parallel$  is in the column space of  $X$ , and  $y^\perp$  is orthogonal to it. Then

$$\|X\beta - y\|_2^2 = \|X\beta - y^\perp - y^\parallel\|_2^2 = \|X\beta - y^\parallel\|_2^2 - \underbrace{2\langle X\beta - y^\perp, y^\parallel \rangle}_0 + \|y^\parallel\|_2^2.$$

Since  $\{X\beta : \beta \in \mathbb{R}^d\}$  is simply the column space of  $X$ , it means there is a choice of  $\beta$  so that  $X\beta = y^\parallel$ , which makes  $\|X\beta - y^\parallel\|_2^2 = 0$ ; this is the optimal choice,  $\beta^{LS}$ . We thus have that  $\beta^{LS}$  is such that  $X\beta^{LS}$  is the orthogonal projection of  $y$  onto the column space of  $X$ , i.e.  $X\beta^{LS} = UU^\top y$  (where  $U$  is from the SVD).

Using the pseudoinverse (see [Definition 6.2.4](#))

$$\beta^{LS} = (X^\top X)^+ X^\top y.$$

since  $X(X^\top X)^+ X^\top$  is the orthogonal projection onto the column space of  $X$  (which can be verified by looking at the SVD).

Computationally, the above most expensive part of computing  $\beta^{LS}$  is computing  $X^\top X$ , which with nested for loops takes  $\Theta(nd^2)$  flops (the pseudoinverse can be computed via the SVD, which itself can be found in  $O(d^3)$  flops). We next show ways of speeding this computation up using sketching.

#### 6.3.1 Sketch-and-solve via subspace embeddings

The *sketch-and-solve* paradigm is a simple one introduced by Sarlós [[Sar06](#)]. The idea is to pick a sketch matrix  $\Pi \in \mathbb{R}^{m \times n}$  ( $m \ll n$ ) and solve for  $\tilde{\beta}^{LS} = \operatorname{argmin} \|\Pi X\beta - \Pi y\|_2^2$  instead.

**Lemma 6.3.1.** *Define  $E := \operatorname{span}(\operatorname{cols}(X), y)$ . Suppose  $\Pi$  is an  $\varepsilon$ -subspace embedding for  $E$ . Then*

$$\|X\tilde{\beta}^{LS} - y\|_2^2 \leq \frac{1 + \varepsilon}{1 - \varepsilon} \cdot \|X\beta^{LS} - y\|_2^2.$$

*Proof.*

$$\begin{aligned} (1 - \varepsilon)\|X\tilde{\beta}^{LS} - y\|_2^2 &\leq \|\Pi X\tilde{\beta}^{LS} - \Pi y\|_2^2 && \text{(subspace embedding property)} \\ &\leq \|\Pi X\beta^{LS} - \Pi y\|_2^2 && (\tilde{\beta}^{LS} \text{ is the minimizer for the sketched problem)} \\ &\leq (1 + \varepsilon)\|X\beta^{LS} - y\|_2^2 && \text{(subspace embedding property)} \end{aligned}$$

The lemma then follows by rearranging terms.  $\square$

Combining with [Subsection 6.2.3](#) for example implies the following theorem of [[CW17](#)].

**Theorem 6.3.2.** *Given  $X \in \mathbb{R}^{n \times d}$ ,  $y \in \mathbb{R}^n$ , in time  $O(\operatorname{nnz}(X) + n) + \operatorname{poly}(d/\varepsilon)$ , with probability at least  $9/10$  one can compute  $\tilde{\beta}^{LS}$  satisfying*

$$\|X\tilde{\beta}^{LS} - y\|_2^2 \leq \frac{1 + \varepsilon}{1 - \varepsilon} \cdot \|X\beta^{LS} - y\|_2^2.$$

Here  $\operatorname{nnz}(A)$  denotes the number of nonzero entries in  $A$ .

*Proof.* Let  $\Pi$  be a CountSketch matrix with  $m = O(d^2/\varepsilon^2)$  rows. Then it is a subspace embedding for  $\text{span}(\text{cols}(X), y)$  with probability at least  $9/10$  by the argument in [Subsection 6.2.3](#). The error guarantee then follows by [Lemma 6.3.1](#). The runtime follows since  $\Pi \cdot$  can be computed in time  $O(m + \text{nnz}(X))$  time, at which point  $(\Pi X)^\top (\Pi X)$  can be computed in time  $O(md^2) = \text{poly}(d/\varepsilon)$ . Then the pseudoinverse can be computed in time  $O(d^3)$  via the SVD. Also,  $\Pi y$  can be computed in time  $O(n)$ . Thus overall  $\tilde{\beta}^{LS} = ((\Pi X)^\top (\Pi X))^\dagger (\Pi X)^\top (\Pi y)$  can be computed in the stated time bound.  $\square$

### 6.3.2 Sketch-and-solve via AMM and subspace embeddings

[Subsection 6.3.1](#) showed how to use  $\varepsilon$ -subspace embeddings to get  $(1 + O(\varepsilon))$ -approximate least squares regression. The downside is then that OSE's would need to have a number of rows depending on  $1/\varepsilon^2$ . In this section we show an alternate analysis of sketch-and-solve due to Sarlós which only requires an  $O(1)$ -subspace embedding (independent of  $\varepsilon$ ), though also requires a certain probabilistic event to hold involving approximate matrix multiplication. In particular, note condition (2) of [Theorem 6.3.3](#) is an AMM guarantee since  $U^\top (X\beta^{LS} - y) = 0$  ( $\beta^{LS}$  is such that  $X\beta^{LS} - y$  is the projection of  $y$  onto the space orthogonal to the column space of  $X$ ).

**Theorem 6.3.3.** *Given the least squares regression problem of minimizing  $\|X\beta - y\|_2^2$  for  $X \in \mathbb{R}^{n \times d}$ , write the SVD  $X = U\Sigma V^\top$ . Suppose  $\Pi$  satisfies the following two conditions:*

- (1)  $\Pi$  is an  $\varepsilon_0$ -subspace embedding for the column space of  $X$ , for  $\varepsilon_0 = 1 - 1/\sqrt{2}$ , and
- (2)  $\|(\Pi U)^\top \Pi(X\beta^{LS} - y)\|_2 \leq \sqrt{\frac{\varepsilon}{2d}} \cdot \|U\|_F \cdot \|X\beta^{LS} - y\|_2$ .

Then if  $\tilde{\beta}^{LS}$  is the minimizer of  $\|\Pi X\beta - \Pi y\|_2^2$ ,

$$\|X\tilde{\beta}^{LS} - y\|_2^2 \leq (1 + O(\varepsilon))\|X\beta^{LS} - y\|_2^2.$$

*Proof.* First some notation: define  $w := y - X\beta^{LS}$ , and  $\alpha, \gamma$  to be vectors such that  $U\alpha = X\beta^{LS}$  and  $U\gamma = X(\tilde{\beta}^{LS} - \beta^{LS})$ . Then

$$\begin{aligned} \|X\tilde{\beta}^{LS} - y\|_2^2 &= \|X\tilde{\beta}^{LS} - X\beta^{LS} + X\beta^{LS} - y\|_2^2 \\ &= \|X(\tilde{\beta}^{LS} - \beta^{LS})\|_2^2 + \|X\beta^{LS} - y\|_2^2 \quad (X\beta^{LS} - y \text{ is orthogonal to } \text{cols}(X)) \\ &= \|\gamma\|_2^2 + \text{OPT}. \end{aligned} \tag{6.3}$$

We would thus like to show that  $\|\gamma\|_2^2 \leq \varepsilon \cdot \text{OPT}$ . Let  $\text{Proj}_A B$  denote the matrix whose columns are the columns of  $B$  projected onto the column space of  $A$ . Observe

$$\begin{aligned} \Pi U(\alpha + \gamma) &= \Pi X\tilde{\beta}^{LS} \\ &= \text{Proj}_{\Pi X}(\Pi y) \\ &= \text{Proj}_{\Pi U}(\Pi y) \quad (\Pi X \text{ and } \Pi U \text{ have the same column space}) \\ &= \text{Proj}_{\Pi U}(\Pi(U\alpha + w)) \\ &= \Pi U\alpha + \text{Proj}_{\Pi U}(\Pi w). \end{aligned} \tag{6.4}$$

Therefore  $\Pi U\gamma = \text{Proj}_{\Pi U}(\Pi w)$  so that  $(\Pi U)^\top (\Pi U)\gamma = (\Pi U)^\top \Pi w$ . Next,  $\Pi$  is an  $\varepsilon_0$ -subspace embedding for  $U$ , we have  $\|(\Pi U)^\top (\Pi U) - I\| \leq \varepsilon_0$  so that the smallest singular value of  $(\Pi U)^\top (\Pi U)$  is at least  $1 - \varepsilon_0$ . Therefore, applying condition (2) in the theorem conditions and using  $\|U\|_F \leq \sqrt{d}$ ,

$$\frac{\|\gamma\|_2^2}{2} = (1 - \varepsilon_0)^2 \|\gamma\|_2^2 \leq \|(\Pi U)^\top \Pi U\gamma\|_2^2 = \|(\Pi U)^\top \Pi w\|_2^2 \leq \frac{\varepsilon}{2} \cdot \text{OPT}. \tag{6.5}$$

Thus  $\|\gamma\|_2^2 \leq \varepsilon \cdot \text{OPT}$ , as desired.  $\square$

### 6.3.3 Accelerating iterative solvers via sketching

In gradient descent we have some function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  for which would like to find an approximate minimizer. The approach is to pick some initial value  $x^{(0)}$  then, iteratively apply the update rule  $x^{(t+1)} \leftarrow x^{(t)} - \gamma \cdot \nabla f(x^{(t)})$ . If one makes assumptions about  $f$ , e.g. we have bounds on the singular values of its Hessian, it is possible to prove a general theorem that shows that  $f(x^{(t)})$  converges exponentially quickly to  $f(x)$ . We will not cover that here, but rather focus specifically on the least squares regression problem.

For least squares regression we define  $f(\beta) = \|X\beta - y\|_2^2$ . As  $f(\beta) = (X\beta - y)^\top (X\beta - y) = \beta^\top X^\top X\beta - 2\beta^\top X^\top y + y^\top y$ , and thus  $\nabla f(\beta) = 2X^\top X\beta - 2X^\top y$ . We will take  $\gamma = 1/2$  and thus have the gradient descent update rule

$$\beta^{(t+1)} = \beta^{(t)} + X^\top (y - X\beta^{(t)}).$$

We now prove a lemma showing that  $f(\beta)$  converges exponentially quickly to  $f(\beta^{LS})$  (recall  $\beta^{LS}$  is the minimizer of  $f$ ) if all singular values of  $X$  are close to 1. This will not be true in general, but as we will soon see, we can make this true quickly using sketching.

**Lemma 6.3.4.** *Suppose  $f(\beta) = \|X\beta - y\|_2^2$  and that all singular values of  $X$  are in the interval  $[1 - 1/\sqrt{2}, 1 + 1/\sqrt{2}]$ . If we perform gradient descent with  $\gamma = 1/2$ , then for all  $t \geq 0$ ,*

$$\|X(\beta^{(t)} - \beta^{LS})\|_2 \leq 2^{-t} \cdot \|X(\beta^{(0)} - \beta^{LS})\|_2. \quad (6.6)$$

*Proof.* The claim holds trivially for  $t = 0$ . For  $t > 0$ , note

$$X(\beta^{(t)} - \beta^{LS}) = X(\beta^{(t-1)} + X^\top (y - X\beta^{(t-1)}) - \beta^{LS}) = (X - XX^\top X)(\beta^{(t-1)} - \beta^{LS}).$$

The final equality holds since  $XX^\top y = XX^\top X\beta^{LS}$  since  $X\beta^{LS}$  is the orthogonal projection of  $y$  onto the column space of  $X$ . Therefore, writing the SVD  $X = U\Sigma V^\top$ ,

$$\begin{aligned} \|X(\beta^{(t)} - \beta^{LS})\|_2 &= \|(X - XX^\top X)(\beta^{(t-1)} - \beta^{LS})\|_2 \\ &= \|U(\Sigma - \Sigma^3)V^\top(\beta^{(t-1)} - \beta^{LS})\|_2 \\ &= \|(I - \Sigma^2)\Sigma V^\top(\beta^{(t-1)} - \beta^{LS})\|_2 && (U \text{ has orthonormal columns}) \\ &\leq \|I - \Sigma^2\| \cdot \|\Sigma V^\top(\beta^{(t-1)} - \beta^{LS})\|_2 && (\|AB\|_F \leq \|A\| \|B\|_F) \\ &= \|I - \Sigma^2\| \cdot \|U\Sigma V^\top(\beta^{(t-1)} - \beta^{LS})\|_2 && (U \text{ has orthonormal columns}) \\ &\leq \frac{1}{2} \cdot \|X(\beta^{(t-1)} - \beta^{LS})\|_2 && (\text{all singular values are } 1 \pm \frac{1}{\sqrt{2}}) \end{aligned}$$

The lemma thus follows by induction on  $t$ . □

Now back to sketching: suppose we want to find some  $\tilde{\beta}$  such that

$$\|X\tilde{\beta} - y\|_2 \leq (1 + \varepsilon)\|X\beta^{LS} - y\|_2.$$

First we will compute  $\beta^{(0)}$  satisfying the above with  $\varepsilon = 1/2$ . This can be done using sketch-and-solve using an  $O(1)$ -subspace embedding. Then note

$$\begin{aligned} \|X(\beta^{(0)} - \beta^{LS})\|_2 &\leq \|X\beta^{(0)} - y\|_2 + \|X\beta^{LS} - y\|_2 && (\text{triangle inequality}) \\ &\leq 2.5\text{OPT} \end{aligned}$$

so that the RHS of Eq. (6.6) is at most  $\varepsilon \cdot \text{OPT}$  for  $t = \Theta(\log(1/\varepsilon))$ . Thus by the triangle inequality,

$$\|X\beta^{(t)} - y\|_2 \leq \|X\beta^{LS} - y\|_2 + \|X(\beta^{(t)} - \beta^{LS})\|_2 \leq (1 + \varepsilon)\text{OPT}$$

for  $t = \Theta(\log(1/\varepsilon))$ . The only thing that remains is to ensure the condition of Lemma 6.3.4 is satisfied: that all the singular values of  $X$  are sufficiently close to 1.

We pick an  $\alpha$ -OSE  $\Pi \in \mathbb{R}^{m \times n}$  for the column space of  $X$  and compute  $\Pi X$ , e.g.  $\Pi$  being the SRHT or the CountSketch so that this computation can be done quickly (see Fig. 6.1). We then compute the SVD  $\Pi X = U'\Sigma'V'^\top$  in time  $O(dm^2)$ . Define  $R := V'\Sigma'^{-1}$ . Then note that for any vector  $z$ ,

$$\|z\|_2^2 = \|U'z\|_2^2 = \|\Pi X Rz\|_2^2 = (1 \pm \alpha)\|XRz\|_2^2.$$

Thus all singular values of  $XR$  are in the interval  $[1/\sqrt{1+\alpha}, 1/\sqrt{1-\alpha}]$ , which is  $1 \pm 1/\sqrt{2}$  for  $\alpha = \Theta(1)$ , thus satisfying the conditions of Lemma 6.3.4. Furthermore  $XR$  and  $X$  have the same column space, and thus we can run gradient descent replacing  $X$  with  $XR$  (then if we find a near-optimal solution  $\tilde{\beta}$ , that corresponds to input  $R\tilde{\beta}$  for the original  $f$ ). The runtime per iteration of gradient descent is only  $O(nnz(X))$  to apply  $X$  and  $X^\top$ , and  $O(d^2)$  to apply  $R$ . Over  $O(\log(1/\varepsilon))$  iterations, this amounts to  $O((nnz(A) + d^2) \log(1/\varepsilon))$  time to run gradient descent, in addition to the time to compute  $\Pi X$  and its SVD.

## 6.4 Approximate low-rank approximation

In the low-rank approximation problem, we are given as input a matrix  $A \in \mathbb{R}^{n \times d}$  and integer  $k > 0$  and would like to compute

$$A_k = \underset{\text{rank}(B) \leq k}{\text{argmin}} \|A - B\|_F.$$

One could ask the question for other norms as well, but in this section we focus exclusively on Frobenius norm. Unlike in the case of regression where we make the assumption  $n \gg d$ , the algorithm we discuss in this section is an improvement even if  $n = d$ ; thus we make no assumption on the relationship between  $n$  and  $d$ . The following is a standard result in linear algebra.

**Theorem 6.4.1** (Eckart-Young theorem). *Write  $A = U\Sigma V^\top$  and let  $U_k$  denote the matrix  $U$  with all but the first  $k$  columns zeroed out (and similarly for  $\Sigma_k, V_k$ ). Then  $A_k = U_k \Sigma_k V_k^\top$ .*

Using Eckart-Young we see that we can compute the best rank- $k$  approximation to  $A$  in  $O(nd^2)$  time, by computing the SVD then truncating. In this section we show a faster approach due to [Sar06], based on sketching. Below, we let  $\text{Proj}_{A,k}(B)$  denote the best rank- $k$  approximation to  $k$  in the column space of  $A$ . We also let  $[B]_k$  denote the best rank- $k$  approximation to  $B$ .

**Theorem 6.4.2.** *For  $A \in \mathbb{R}^{n \times d}$  with SVD  $A = U\Sigma V^\top$  and  $k > 0$  an integer, let  $\Pi$  satisfy the following two conditions:*

- (1)  $\Pi$  is an  $\varepsilon_0$ -subspace embedding for the column space of  $V_k$  for  $\varepsilon_0 = 1 - 1/\sqrt{2}$ , and
- (2)  $\|(\Pi V_k)^\top \Pi(A - A_k)\|_F \leq \sqrt{\frac{\varepsilon}{2k}} \cdot \|V_k\|_F \cdot \|A - A_k\|_F$ .

Then if  $\tilde{A}_k := \text{Proj}_{A\Pi^\top, k}(A)$ ,

$$\|A - \tilde{A}_k\|_F^2 \leq (1 + \varepsilon)\|A - A_k\|_F^2.$$

*Proof.* Let  $P$  be the orthogonal projection onto the column space of  $\text{Proj}_{A\Pi^\top}(A_k)$  so that  $P = (\text{Proj}_{A\Pi^\top}(A_k))(\text{Proj}_{A\Pi^\top}(A_k))^\top$ . Since  $\text{Proj}_{A\Pi^\top, k}(A)$  is the best rank- $k$  approximation to  $A$  in the column space of  $A\Pi^\top$  and  $PA$  is in this space,

$$\|A - \text{Proj}_{A\Pi^\top, k}(A)\|_F^2 \leq \|A - PA\|_F^2,$$

so we bound the RHS.

$$\begin{aligned} \|A - PA\|_F^2 &= \|(I - P)(A - A_k) + (I - P)A_k\|_F^2 \\ &= \|(I - P)(A - A_k)\|_F^2 + \|(I - P)A_k\|_F^2 \\ &\leq \|I - P\| \cdot \|A - A_k\|_F^2 + \|(I - P)A_k\|_F^2 \quad (\|AB\|_F \leq \|A\| \cdot \|B\|_F) \\ &= \|A - A_k\|_F^2 + \|(I - P)A_k\|_F^2 \end{aligned}$$

It thus suffices to show  $\|(I - P)A_k\|_F^2 \leq \varepsilon\|A - A_k\|_F^2$ . By definition of  $P$ , we have  $PA_k = (A\Pi^\top)(A\Pi^\top)^\top A_k$  is the best approximation of  $A_k$  in the column space of  $A\Pi^\top$  under Frobenius norm error. Thus, letting  $M^{(i)}$  denote the  $i$ th column of a matrix  $M$ ,

$$\begin{aligned} \|A_k - PA_k\|_F^2 &= \|A_k - (A\Pi^\top)^\top (A\Pi^\top)^\top A_k\|_F^2 \\ &\leq \|A_k - (A\Pi^\top)^\top (A_k\Pi^\top)^\top A_k\|_F^2 \\ &= \|(A\Pi^\top)^\top (A_k\Pi^\top)^\top A_k - A_k\|_F^2 \\ &= \|A_k^\top (\Pi A_k^\top)^\top (\Pi A^\top) - A_k^\top\|_F^2 \\ &= \sum_{i=1}^n \|A_k^\top (\Pi A_k^\top)^\top (\Pi A^\top)^{(i)} - A_k^{\top(i)}\|_2^2 \end{aligned} \quad (6.7)$$

The above is related to  $n$  regression problems. Specifically, consider the regression problems

$$X^{*(i)} = \underset{X}{\text{argmin}} \|A_k^\top X^{(i)} - A^{\top(i)}\|_2^2 = A_k^{\top(i)}$$

and the sketch-and-solve optima

$$\tilde{X}^{*(i)} = \underset{X}{\text{argmin}} \|\Pi A_k^\top X^{(i)} - \Pi A^{\top(i)}\|_2^2$$

for  $i = 1, 2, \dots, n$ . Then  $\tilde{X}^{*(i)} = (\Pi A_k^\top)^\top (\Pi A^\top)^{(i)}$ . Thus [Eq. \(6.7\)](#) can be rewritten as

$$\sum_{i=1}^n \|A_k^\top (\tilde{X}^{*(i)} - X^{*(i)})\|_2^2 = \|A_k^\top (\tilde{X}^* - X^*)\|_F^2,$$

treating  $X^*$  as a matrix with columns  $X^{*(i)}$  (and similarly for  $\tilde{X}^*$ ). The completion of the proof is then essentially identical to that of [Theorem 6.3.3](#). Specifically, note  $A_k^\top = V_k \Sigma_k U_k^\top$ . Thus we can write  $X^* = V_k Z$  for some matrix  $Z$  (its columns are in the column space of  $V_k$ ) and  $A_k^\top (\tilde{X}^* - X^*) = V_k \Gamma$  for some matrix  $\Gamma$ . Also define  $W := A^\top - A^\top X^*$ . These definitions parallel those in the proof of [Theorem 6.3.3](#); specifically  $W$  plays the role of  $w$ ,  $Z$  of  $\alpha$ , and  $\Gamma$  of  $\gamma$ . Then just as in [Eq. \(6.3\)](#), we must show  $\|\Gamma\|_F^2 \leq \varepsilon \|A_k^\top X^* - A^\top\|_F^2 = \varepsilon \|A - A_k\|_F^2$ . The exact same line of reasoning as [Eq. \(6.4\)](#) then implies  $(\Pi V_k)^\top \Pi V_k \Gamma = (\Pi V_k)^\top \Pi W$ , and using that  $\Pi$  is a subspace embedding for  $V_k$ , similarly as in [Eq. \(6.5\)](#) and using (2) together with the fact that  $\|V_k\|_F = \sqrt{k}$ ,

$$\frac{\|\Gamma\|_F^2}{2} = (1 - \epsilon_0)^2 \|\Gamma\|_F^2 \leq \|(\Pi V_k)^\top \Pi V_k \Gamma\|_F^2 = \|(\Pi V_k)^\top \Pi W\|_F^2 \leq \frac{\varepsilon}{2} \cdot \|A - A_k\|_F^2.$$

□



Now in terms of computation time, we would like to compute  $\text{Proj}_{A\Pi^\top, k}(A)$ . For  $\Pi \in \mathbb{R}^{m \times d}$ , we can compute  $A\Pi^\top$  quickly if  $\Pi$  is either the SRHT or CountSketch (see Fig. 6.1); note for these choices of matrices and the conditions in Theorem 6.4.2, we can take  $m = \tilde{O}(k/\varepsilon)$  (SRHT) or  $O(k^2 + k/\varepsilon)$  (CountSketch). Then  $A\Pi^\top$  is an  $n \times m$  matrix, and we can thus compute its SVD  $U'\Sigma'V'^\top$  in time  $O(nm^2)$ . We then wish to compute  $[U'U'^\top A]_k = U'[U'^\top A]_k$ . We can compute  $U'^\top A$  in  $O(nmd)$  time (since  $U'$  has rank at most  $m$  and thus has at most  $m$  columns), then computing its SVD takes time  $O(dm^2)$ . The overall runtime is faster than the full  $O(nd^2)$  time to compute the SVD, as we replace one factor of  $d$  with  $m$ , which is on the order of  $k$  or  $k^2$  (note if  $d > n$ , we can apply the algorithm of this section to  $A^\top$ ).

## 6.5 Projection-cost preserving sketches

Another method for sketching for low-rank approximation, introduced by Cohen, Elder, Musco, Musco, and Persu [CEM<sup>+</sup>15], is that of using a *projection-cost preserving (PCP) sketch*.

**Definition 6.5.1.** Let  $\mathcal{O}_k$  be the set of all orthogonal projections onto subspaces of dimension at most  $k$ . Given a matrix  $A \in \mathbb{R}^{n \times d}$ , we say  $\Pi \in \mathbb{R}^{m \times d}$  is an  $(\varepsilon, k)$ -PCP for  $A$  if

$$\forall P \in \mathcal{O}_k, (1 - \varepsilon)\|(I - P)A\|_F^2 \leq \|(I - P)A\Pi^\top\|_F^2 \leq (1 + \varepsilon)\|(I - P)A\|_F^2.$$

Using a PCP sketch allows for a “sketch-and-solve” type of approach for low-rank approximation, which is different from the algorithm described in Section 6.4. That is, suppose  $\mathcal{R}_k \subseteq \mathcal{O}_k$  and one wishes to find

$$P^* = \underset{P \in \mathcal{R}_k}{\operatorname{argmin}} \|(I - P)A\|_F^2.$$

One can then simply instead compute

$$\tilde{\mathcal{P}}^* = \underset{P \in \mathcal{R}_k}{\operatorname{argmin}} \|(I - P)A\Pi^\top\|_F^2,$$

and an analysis similar to the proof of Lemma 6.3.1 implies

$$\|(I - \tilde{\mathcal{P}}^*)A\|_F^2 \leq \frac{1 + \varepsilon}{1 - \varepsilon} \|(I - P^*)A\|_F^2.$$

If  $\mathcal{R}_k = \mathcal{O}_k$ , this gives a sketch-and-solve solution to the same problem studied in Section 6.4 (approximate unconstrained low-rank approximation under the Frobenius norm). However, other problems can also be captured via different  $\mathcal{R}_k$ , showing that a sketch-and-solve approach works for them as well. Consider for example the *k-means clustering problem* mentioned in Chapter 5. In this problem we are given an integer parameter  $k > 0$  as well as  $n$  points  $x_1, \dots, x_n \in \mathbb{R}^d$ . We must then find  $k$  cluster centers  $y_1, \dots, y_k \in \mathbb{R}^d$  minimizing

$$\sum_{i=1}^n \min_{1 \leq j \leq k} \|x_i - y_j\|_2^2.$$

As discussed in Eq. (5.2), the  $k$ -means objective can be rewritten as computing

$$\mathcal{P}^* = \underset{k\text{-partitions } \mathcal{P} \text{ of } [n]}{\operatorname{argmin}} \sum_{j=1}^k \sum_{i \in \mathcal{P}_j} \|x_i - \mu_j\|_2^2$$

for  $\mu_j := (1/|\mathcal{P}_j|) \sum_{i \in \mathcal{P}_j} x_i$ , where  $\mathcal{P} = (\mathcal{P}_1, \dots, \mathcal{P}_k)$ . An observation of [BZMD15] is that if one defines  $X_{\mathcal{P}}$  as the  $n \times k$  matrix with

$$(X_{\mathcal{P}})_{i,j} = \begin{cases} \frac{1}{\sqrt{|\mathcal{P}_j|}} & i \in \mathcal{P}_j \\ 0 & \text{otherwise} \end{cases},$$

then the columns of are orthonormal, so that  $Z_{\mathcal{P}} := X_{\mathcal{P}}X_{\mathcal{P}}^{\top}$  is a rank- $k$  orthogonal projection<sup>1</sup>. Also, if one puts the  $x_i$  as rows of a matrix  $A$ , then a computation shows that the  $i$ th row of  $Z_{\mathcal{P}}A$  equals  $\mu_j$ , where  $i \in \mathcal{P}_j$ . Thus if we define  $\mathcal{Q}_k$  as the set of all  $Z_{\mathcal{P}}$  for  $\mathcal{P}$  a  $k$ -partition, then the  $k$ -means objective can be rewritten as computing

$$Z^* = \underset{Z \in \mathcal{Q}_k}{\operatorname{argmin}} \|(I - Z)A\|_F^2,$$

which can be solved to approximate optimality using sketch-and-solve with a PCP sketch. As we will soon see, oblivious randomized PCP sketches exist with  $O(k/\varepsilon^2)$  rows, meaning that  $k$ -means can be approximately solved by oblivious projection down to  $O(k/\varepsilon^2)$  dimensions, independent of  $n$  (unlike the discussion surrounding Eq. (5.3), which showed that dimension  $O(\varepsilon^{-2} \log n)$  suffices). In fact it has been shown even more recently [MMR19], via an argument unrelated to this section, that projecting to dimension  $O(\log(k/\varepsilon)/\varepsilon^2)$  in fact suffices to preserve the cost of any  $k$ -clustering of a set of points, independent of  $n$ .

We now show the main theorem of this section. If we choose  $\Pi \in \mathbb{R}^{m \times d}$  obliviously to  $A$  from some distribution  $\mathcal{D}$ , we state next to each of the bulleted conditions in the theorem statement what  $\mathcal{D}$  should satisfy so all conditions are simultaneously satisfied with probability at least  $1 - \delta$ .

**Theorem 6.5.2** ([CEM<sup>+</sup>15]). *Let  $A \in \mathbb{R}^{n \times d}$  be given, and suppose its SVD decomposition is  $A = U\Sigma V^{\top}$ . Let  $A_{-k}$  denote  $A - A_k = A - U_k \Sigma_k V_k^{\top}$ . If  $\Pi \in \mathbb{R}^{m \times d}$  satisfies the following four properties, then it is a  $(5\varepsilon, k)$ -PCP sketch for  $A$ :*

- $\Pi$  is an  $\varepsilon$ -subspace embedding for the rowspace of  $A_k$ . *Implied by  $(\varepsilon, k, \frac{\delta}{4})$ -OSE.*
- $|\|A_{-k}\Pi^{\top}\|_F^2 - \|A_{-k}\|_F^2| \leq \varepsilon \|A_{-k}\|_F^2$ .  *$(\varepsilon, \frac{\delta}{4}, p)$ -JLMP for  $p \geq 1$ ; Lemma 6.5.3.*
- $\|A_{-k}A_{-k}^{\top} - A_{-k}\Pi^{\top}\Pi A_{-k}^{\top}\|_F \leq \frac{\varepsilon}{\sqrt{k}} \cdot \|A_{-k}\|_F^2$ .  *$(\frac{\varepsilon}{\sqrt{k}}, \frac{\delta}{4}, p)$ -JLMP for  $p \geq 2$ ; Theorem 6.1.4.*
- $\|(\Pi V_k)^{\top}(\Pi A_{-k}^{\top})\|_F \leq \frac{\varepsilon}{\sqrt{k}} \cdot \|V_k\|_F \cdot \|A_{-k}\|_F$ .  *$(\frac{\varepsilon}{\sqrt{k}}, \frac{\delta}{4}, p)$ -JLMP for  $p \geq 2$ ; Theorem 6.1.4.*

*Proof.* Our goal is to show that under the conditions of the theorem statement,

$$\forall P \in \mathcal{O}_k, (1 - \varepsilon)\|(I - P)A\|_F^2 \leq \|(I - P)A\Pi^{\top}\|_F^2 \leq (1 + \varepsilon)\|(I - P)A\|_F^2.$$

Alternatively by writing  $Y = I - P$ , it is equivalent to say

$$\forall Y \in \mathcal{O}_{n-k}, (1 - \varepsilon) \operatorname{tr}(YAA^{\top}Y) \leq \operatorname{tr}(YA\Pi^{\top}\Pi A^{\top}Y) \leq (1 + \varepsilon) \operatorname{tr}(YAA^{\top}Y).$$

We expand  $A = A_k + A_{-k}$ . Then

$$\operatorname{tr}(YAA^{\top}Y) = \operatorname{tr}(YA_kA_k^{\top}Y) + \operatorname{tr}(YA_{-k}A_{-k}^{\top}Y) + \underbrace{\operatorname{tr}(YA_{-k}A_k^{\top}Y)}_0 + \underbrace{\operatorname{tr}(YA_kA_{-k}^{\top}Y)}_0$$

<sup>1</sup>We eliminate empty  $P_j$ , so we may have a  $k'$ -partition for some  $k' < k$ , and thus a rank- $k'$  projection.

and

$$\begin{aligned} \text{tr}(Y A \Pi^\top \Pi A^\top Y) &= \text{tr}(Y A_k \Pi^\top \Pi A_k^\top Y) + \text{tr}(Y A_{-k} \Pi^\top \Pi A_{-k}^\top Y) \\ &\quad + \text{tr}(Y A_{-k} \Pi^\top \Pi A_k^\top Y) + \text{tr}(Y A_k \Pi^\top \Pi A_{-k}^\top Y). \end{aligned}$$

Thus if we define

- $E_1 := A_k A_k^\top - A_k \Pi^\top \Pi A_k^\top$
- $E_2 := A_{-k} A_{-k}^\top - A_{-k} \Pi^\top \Pi A_{-k}^\top$
- $E_3 := A_{-k} \Pi^\top \Pi A_k^\top$
- $E_4 := A_k \Pi^\top \Pi A_{-k}^\top$ ,

we would like to show  $|\text{tr}(Y E_i Y)| \leq \varepsilon \cdot \text{tr}(Y C Y)$  for each  $i = 1, \dots, 4$  for  $C = A A^\top$ . Note also that  $\|A_{-k}\|_F^2 \leq \text{tr}(Y C Y)$ , as  $A_{-k} = (I - U_k U_k^\top) A$  where  $P^*$  is the minimizer of  $\|(I - P)A\|_F^2$  over all rank- $k$  orthogonal projections.

**E<sub>1</sub>** : Define  $y_i$  to be the  $i$ th column of  $Y$ . Then

$$\begin{aligned} |\text{tr}(Y E_1 Y)| &= \left| \sum_i \|A_k^\top y_i\|_2^2 - \|\Pi A_k^\top y_i\|_2^2 \right| \\ &\leq \varepsilon \cdot \sum_i \|A_k^\top y_i\|_2^2 \quad (\Pi \text{ is an } \varepsilon\text{-subspace embedding for rowspace of } A_k) \\ &= \varepsilon \cdot \text{tr}(Y C_k Y) \\ &\leq \varepsilon \cdot \text{tr}(Y C Y). \end{aligned}$$

**E<sub>2</sub>** :

$$\begin{aligned} |\text{tr}(Y E_2 Y)| &= |\text{tr}(Y E_2)| \\ &= |\text{tr}(E_2) - \text{tr}(P E_2)| \\ &\leq |\text{tr}(E_2)| + |\text{tr}(P E_2)| \end{aligned}$$

We have

$$|\text{tr}(E_2)| = \|\|A_{-k} \Pi^\top\|_F^2 - \|A_{-k}\|_F^2\| \leq \varepsilon \|A_{-k}\|_F^2 \leq \varepsilon \text{tr}(Y C Y).$$

Observe  $\text{tr}(P E_2)$  is the sum of eigenvalues of a rank at most  $k$  matrix. Hence by Cauchy-Schwarz,

$$|\text{tr}(P E_2)| \leq \sqrt{k} \cdot \|P E_2\|_F \leq \sqrt{k} \cdot \|E_2\|_F \leq \varepsilon \|A_{-k}\|_F^2 \leq \varepsilon \text{tr}(Y C Y).$$

**E<sub>3</sub>** : Since  $E_3^\top$  is in the column space of  $A$ ,  $C C^+ E_3^\top = E_3^\top$ . Thus

$$\begin{aligned} \text{tr}(Y E_3 Y) &= \text{tr}(Y E_3^\top Y) \\ &= \text{tr}(Y E_3^\top) \quad (\text{tr}(AB) = \text{tr}(BA), \text{ and } Y Y = Y \text{ since } Y \text{ is a projection}) \\ &= \text{tr}(Y C C^+ E_3^\top) \\ &= \text{tr}(Y C (C^+)^{1/2} (C^+)^{1/2} E_3^\top) \tag{6.8} \\ &\leq \sqrt{\text{tr}(Y C C^+ C Y)} \cdot \sqrt{\text{tr}(E_3 C^+ E_3^\top)} \quad (\text{Cauchy-Schwarz}) \end{aligned}$$

$$= \sqrt{\text{tr}(YCY)} \cdot \sqrt{\text{tr}(E_3 C^+ E_3^\top)} \quad (6.9)$$

In Eq. (6.8) we used that  $C = U\Sigma^2 U^\top$ , so  $C^+ = U\Sigma^{-2} U^\top$  and thus  $(C^+)^{1/2} = U\Sigma^{-1} U^\top$ . Also,

$$E_3 = U\Sigma V_{-k}^\top \Pi^\top \Pi V_k \Sigma U^\top.$$

Then

$$\begin{aligned} \text{tr}(E_3 C^+ E_3^\top) &= \text{tr}(U\Sigma V_{-k}^\top \Pi^\top \Pi V_k \Sigma U^\top U\Sigma^{-2} U^\top U\Sigma V_k^\top \Pi^\top \Pi V_{-k} \Sigma U^\top) \\ &= \text{tr}(A_{-k} \Pi^\top \Pi V_k V_k^\top \Pi^\top \Pi A_{-k}^\top) \\ &= \|(\Pi V_k)^\top (\Pi A_{-k}^\top)\|_F^2 \\ &\leq \frac{\varepsilon^2}{k} \cdot \|V_k\|_F^2 \cdot \|A_{-k}\|_F^2 \\ &= \varepsilon^2 \cdot \|A_{-k}\|_F^2. \end{aligned}$$

Thus

$$\text{Eq. (6.9)} \leq \varepsilon \cdot \sqrt{\text{tr}(YCY)} \cdot \|A_{-k}\|_F \leq \varepsilon \cdot \text{tr}(YCY)$$

**E<sub>4</sub>** :  $E_4 = E_3^\top$ , so  $|\text{tr}(Y E_4 Y)| = |\text{tr}(Y E_3 Y)| \leq \varepsilon \cdot \text{tr}(YCY)$ . □

The four conditions of [Theorem 6.5.2](#) should look familiar. The first is a simple subspace embedding requirement. The third and fourth conditions are AMM conditions under Frobenius norm with error  $\varepsilon/\sqrt{k}$ , which happen with good probability if  $\Pi$  satisfies the  $(\varepsilon/\sqrt{k}, \delta/4, p)$ -JL moment property for some  $p \geq 2$  (see [Theorem 6.1.4](#)). The second bullet looks quite a bit like what would be implied by the Distributional JL lemma, although applied not to a vector but to the matrix  $A_{-k}^\top$  (note  $\|A_{-k} \Pi^\top\|_F^2 = \|\Pi A_{-k}^\top\|_F^2$ ). We show in [Lemma 6.5.3](#) that the second bullet holds with good probability when  $\Pi$  is drawn from a distribution with the  $(\varepsilon, \delta/4, p)$ -JL moment property for some  $p \geq 1$ . Thus in summary: we would like  $\Pi$  to be drawn from a distribution that is simultaneously an  $(\varepsilon, k)$ -OSE, and also has the  $(\varepsilon/\sqrt{k}, \delta/4, p)$ -JL moment property for some  $p \geq 2$ . The CountSketch with  $m$  rows satisfies these conditions with  $m = O(d^2/(\varepsilon^2 \delta))$ , allowing for fast multiplication, for example, or one could also use the OSNAP or SRHT distributions, or a matrix  $\Pi$  with i.i.d. gaussian entries (see [Fig. 6.1](#)).

**Lemma 6.5.3.** *Suppose  $\Pi \in \mathbb{R}^{m \times n}$  is drawn from a distribution satisfying the  $(\varepsilon, \delta, p)$ -JL moment property for some  $p \geq 1$ . Then for any matrix  $M \in \mathbb{R}^{n \times d}$ ,*

$$\mathbb{P}_{\Pi}(|\|\Pi M\|_F^2 - \|M\|_F^2| > \varepsilon \|M\|_F^2) < \delta.$$

*Proof.* Let  $m_i$  denote the  $i$ th column of  $M$ . Then  $\|\Pi M\|_F^2 = \sum_i \|\Pi m_i\|_2^2$ . Thus

$$\begin{aligned} \|\|\Pi M\|_F^2 - \|M\|_F^2\|_p &= \left\| \sum_i (\|\Pi m_i\|_2^2 - \|m_i\|_2^2) \right\|_p \\ &\leq \sum_i \|\|\Pi m_i\|_2^2 - \|m_i\|_2^2\|_p && \text{(triangle inequality)} \\ &\leq \sum_i \|m_i\|_2^2 \cdot \varepsilon \delta^{1/p} && \text{(JL moment property)} \\ &= \varepsilon \delta^{1/p} \cdot \|M\|_F^2. \end{aligned}$$

Thus by Markov's inequality,

$$\mathbb{P}(|\|\Pi M\|_F^2 - \|M\|_F^2| > \varepsilon \|M\|_F^2) < \frac{\|\|\Pi M\|_F^2 - \|M\|_F^2\|_p^p}{\varepsilon^p \|M\|_F^{2p}} \leq \delta.$$

□

### 6.5.1 $k$ -means clustering

As mentioned already, a recent result shows that in fact one can do sketch-and-solve for  $k$  means when randomly projecting down to dimension  $O(\log(k/\varepsilon)/\varepsilon^2)$  [MMR19], e.g. via a random gaussian matrix. Their result also extends to  $k$ -median, where one sums the distances to cluster centers and not the squared distances. We will not cover their result here, but we will show a precursor result that randomly projecting to  $O(\varepsilon^{-2} \log k)$  dimensions implies that sketch-and-solve will find a clustering that is approximately optimal up to a factor  $9 + \varepsilon$  [CEM<sup>+</sup>15] (instead of  $1 + \varepsilon$  as in [MMR19]).

Recall in  $k$ -means we are given a matrix  $A \in \mathbb{R}^{n \times d}$ , whose rows we would like to cluster, and the goal is to find some  $P$  in  $\mathcal{Q}_k$  which minimizes  $\|(I - P)A\|_F^2$ .

In the below lemma, we again in blue write what conditions the distribution  $\Pi$  is drawn from to ensure that all conditions are simultaneously satisfied with probability at least  $1 - \delta$ . For the second condition, note that the rows of  $P^*A$  are the centroids of the corresponding points after clustering using  $P^*$ , and thus  $P^*A$  only has  $k$  distinct rows. One can see that a gaussian sketching matrix  $\Pi$ , for example, satisfies the desired properties with  $m = O(\log(k/\delta)/\varepsilon^2)$  rows.

**Lemma 6.5.4** ([CEM<sup>+</sup>15]). *Let  $A \in \mathbb{R}^{n \times d}$  be given. Suppose  $\Pi \in \mathbb{R}^{m \times d}$ . Let  $P^*$  denote the optimal  $k$ -means clustering projection matrix for  $A$ , and  $\tilde{P}^*$  for  $A\Pi^\top$ . Let  $\tilde{P} \in \mathbb{R}^{n \times n}$  be a clustering projection matrix. Suppose  $\Pi \in \mathbb{R}^{m \times d}$  satisfies the following conditions:*

- $\|(I - \tilde{P})P^*A\|_F^2 \leq (1 + \varepsilon)\|(I - \tilde{P})P^*A\Pi^\top\|_F^2$ . ( $(\frac{\varepsilon}{2}, \frac{\delta}{3k^2}, p)$ -JLMP for  $p \geq 1$ ; Eq. (5.3).
- $\|(I - \tilde{P}^*)A\Pi^\top\|_F^2 \leq (1 + \varepsilon)\|(I - \tilde{P}^*)A\|_F^2$ . ( $(\varepsilon, \frac{\delta}{3}, p)$ -JLMP for  $p \geq 1$ ; Lemma 6.5.3.
- $\|(I - P^*)A\Pi^\top\|_F^2 \leq (1 + \varepsilon)\|(I - P^*)A\|_F^2$ . ( $(\varepsilon, \frac{\delta}{3}, p)$ -JLMP for  $p \geq 1$ ; Lemma 6.5.3.

Then if  $\tilde{P}$  satisfies

$$\|(I - \tilde{P})A\Pi^\top\|_F^2 \leq \gamma\|(I - \tilde{P}^*)A\Pi^\top\|_F^2,$$

i.e. it is a  $\gamma$ -approximately optimal  $k$ -means clustering solution for the sketched input  $A\Pi^\top$ , then

$$\|(I - \tilde{P})A\|_F^2 \leq (9 + 12\varepsilon + 4\varepsilon^2) \cdot \gamma\|(I - P^*)A\|_F^2.$$

*Proof.* Write  $B = P^*A$ ,  $\bar{B} = A - B = (I - P^*)A$ . Then

$$\begin{aligned} \|A - \tilde{P}A\|_F &= \|(I - \tilde{P})B + (I - \tilde{P})\bar{B}\|_F \\ &\leq \|(I - \tilde{P})B\|_F + \|(I - \tilde{P})\bar{B}\|_F && \text{(triangle inequality)} \\ &\leq \|(I - \tilde{P})B\|_F + \|\bar{B}\|_F && ((I - \tilde{P}) \text{ is a projection}) \\ &\leq \sqrt{1 + \varepsilon}\|(I - \tilde{P})B\Pi^\top\|_F + \|\bar{B}\|_F && \text{(first condition in lemma statement)} \\ &= \sqrt{1 + \varepsilon}\|(A\Pi^\top - \bar{B}\Pi^\top) - \tilde{P}(A\Pi^\top - \bar{B}\Pi^\top)\|_F + \|\bar{B}\|_F && (B\Pi^\top = A\Pi^\top - \bar{B}\Pi^\top) \\ &\leq \sqrt{1 + \varepsilon}\|(I - \tilde{P})A\Pi^\top\|_F + \sqrt{1 + \varepsilon}\|(I - \tilde{P})\bar{B}\Pi^\top\|_F + \|\bar{B}\|_F && \text{(triangle inequality)} \\ &\leq \sqrt{1 + \varepsilon}\|(I - \tilde{P})A\Pi^\top\|_F + \sqrt{1 + \varepsilon}\|\bar{B}\Pi^\top\|_F + \|\bar{B}\|_F && (I - \tilde{P} \text{ is a projection matrix}) \end{aligned}$$

$$\begin{aligned}
&\leq \sqrt{1+\varepsilon}\sqrt{\gamma}\|(I-\tilde{P}^*)A\Pi^\top\|_F + \sqrt{1+\varepsilon}\|\bar{B}\Pi^\top\|_F + \|\bar{B}\|_F \quad (\tilde{P}^* \text{ is optimal for } A\Pi^\top) \\
&\leq (1+\varepsilon)\sqrt{\gamma}\|(I-\tilde{P}^*)A\|_F + (1+\varepsilon)\|\bar{B}\|_F + \|\bar{B}\|_F \quad (\text{conditions of lemma statement}) \\
&\leq (1+\varepsilon)\sqrt{\gamma}\|(I-P^*)A\|_F + (1+\varepsilon)\|\bar{B}\|_F + \|\bar{B}\|_F \quad (P^* \text{ is optimal for } A) \\
&\leq (3+2\varepsilon)\sqrt{\gamma}\|(I-P^*)A\|_F \quad (\bar{B} = (I-P^*)A, \text{ and } \gamma \geq 1)
\end{aligned}$$

Thus  $\|A - \tilde{P}A\|_F^2 \leq (3+2\varepsilon)^2\gamma\|(I-P^*)A\|_F^2$ , as desired.  $\square$

## Chapter 7

# Compressed Sensing

The field of *compressed sensing* is concerned with (approximately) recovering a signal that is (approximately) sparse in some known basis, based on a sublinear number of measurements. In other words, if we arrange these measurements as rows of a matrix  $\Pi \in \mathbb{R}^{m \times n}$ , then given  $\Pi x$  for some  $x \in \mathbb{R}^n$  ( $m \ll n$ ) we would like recover some  $\tilde{x}$  given only access to  $\Pi x$  and  $\Pi$  such that  $\|x - \tilde{x}\|$  is small, where this smallness depends on how close  $x$  is to being sparse. Recovery guarantees sought are typically of the form

$$\|x - \tilde{x}\|_A \leq f(k) \cdot \min_{\|z\|_0 \leq k} \|x - z\|_B$$

for some norms  $\|\cdot\|_A, \|\cdot\|_B$  where  $k$  is the sparsity parameter, so that  $z$  is the best  $k$ -sparse approximation to  $x$  under the  $\|\cdot\|_B$  norm ( $\|\cdot\|_0$  denotes support size). For example if  $x$  actually is  $k$ -sparse, then the right hand side is zero (simply choose  $z = x$ ). In the norms we consider this chapter ( $\ell_p$  norms), the optimal  $z$  is simply the projection of  $x$  onto its top  $k$  coordinates in magnitude, i.e.  $x_{head(k)}$  (which is  $x - x_{tail(k)}$ , to use the definition of [Remark 4.1.1](#)). Thus the right hand side becomes  $f(k) \cdot \|x_{tail(k)}\|_p$  for some  $p$ .

A good example of approximate sparsity is that of images. An  $n \times n$  pixelated image can be seen as a  $3n^2$  dimensional vector, where each of the  $n^2$  pixels corresponds to three dimensions (RGB values). The color black is represented by the RGB tuple (0,0,0). Now, most interesting images are not simply all black, i.e. they are not sparse (or even approximately sparse) in the *standard* basis. However, most natural images are in fact sparse in the 2D Haar Wavelet basis. Rather than describe this basis in terms of its basis vectors, we instead describe how to perform a change into this basis from the standard basis. The main intuition is that in natural images, typically adjacent pixels are of the same or similar color, except potentially for boundaries between objects (but most pixels are not at a boundary). The change of basis into the 2D Haar Wavelet basis for  $n^2$ -dimensional vectors (which we view as  $n \times n$  images) is thus as follows: we assume  $n$  is a power of 2. Imagine the image is grayscale so that each pixel just has one associated value instead of three values (RGB); the case of color images can be handled by applying the following transform to each color separately. We describe how to transform the input image into the output image after applying the transform. Divide the image pixels into  $(n/2)^2$  blocks, each block of which is  $2 \times 2$ . The output image can be viewed as four images (top left, top right, bottom left, and bottom right), each of which is  $(n/2) \times (n/2)$ , where each input block gives rise to four pixels: one pixel per sub-image in the output block. Specifically if an input block has pixel values  $p_1, p_2, p_3, p_4$ , then the top left image in the output has pixel value equal to their average  $(p_1 + p_2 + p_3 + p_4)/4$ . The other three output sub-images have pixel values  $(p_1 + p_2 - p_3 - p_4)/4$ ,  $(p_1 - p_2 + p_3 - p_4)/4$ , and  $(p_1 - p_2 - p_3 + p_4)/4$ , respectively. We then recursively apply the transform to the top left sub-image of the output. The base case is when  $n = 1$ , in which case we output the 1 input pixel itself. This

procedure is invertible, and thus the linear transformation we describe has full rank, though in practice since pixel values must be integers in  $[0, 255]$ , actual implementations are not invertible due to rounding when dividing by 4.



Figure 7.1: The original image was 2048x2048 pixels. Image (a) applies just one level of the 2D Haar Wavelet Transform on each of the RGB values per pixel separately, and (b) applies the full recursion down to the 1x1 base case. The resulting image is approximately sparse (RGB values of (0,0,0) correspond to black). If one zooms into this document and looks closely, some object edges are still visible after the transform (the outline of my head and jacket collar), albeit faintly.

---

```

from PIL import Image

def haar(pixels, n):
    if n == 1:
        return
    else:
        L = [ [None]*n for _ in range(n) ]
        for i in xrange(n):
            for j in xrange(n):
                L[i][j] = pixels[i,j]

        for i in xrange(n/2):
            for j in xrange(n/2):
                pixels[i,j] = tuple((L[2*i][2*j][k]+L[2*i][2*j+1][k]+L[2*i+1][2*j+1][k]+L[2*i+1][2*j][k])/4 for k in range(3))
                pixels[i,j+n/2] = tuple((L[2*i][2*j][k]+L[2*i][2*j+1][k]-L[2*i+1][2*j+1][k]-L[2*i+1][2*j][k])/4 for k in range(3))
                pixels[i+n/2,j] = tuple((L[2*i][2*j][k]-L[2*i][2*j+1][k]+L[2*i+1][2*j+1][k]-L[2*i+1][2*j][k])/4 for k in range(3))
                pixels[i+n/2,j+n/2] = tuple((L[2*i][2*j][k]-L[2*i][2*j+1][k]-L[2*i+1][2*j+1][k]+L[2*i+1][2*j][k])/4 for k in range(3))

        haar(pixels, n/2)

im = Image.open('input.jpg')
n = im.size[0] # we assume this is a power of 2
haar(im.load(), n)
im.save('output.jpg')

```

---

Figure 7.2: Python code used to generate image (b) in Fig. 7.1, using the Python Imaging Library.

The types of schemes developed in the compressed sensing literature generally fall into one of two categories: *nonuniform* (also known as *for-each*) schemes, and *uniform* (also known as *for-all*) schemes. If we let  $\mathcal{A}$  denote the recovery algorithm which returns  $\tilde{x}$  from  $\Pi x$  (imagine that  $\mathcal{A}$  has  $\Pi$  hardcoded into its source code), then the difference between these two categories is as follows:

- **Nonuniform:** These are randomized schemes which satisfy the guarantee

$$\forall x \in \mathbb{R}^n, \mathbb{P}_{\Pi, \mathcal{A}} (\mathcal{A}(\Pi x) \text{ is a good approximation to } x) \geq 1 - \delta.$$

- **Uniform:** These schemes are either deterministic, or if randomized satisfy the guarantee

$$\mathbb{P}_{\Pi, \mathcal{A}} (\forall x \in \mathbb{R}^n \mathcal{A}(\Pi x) \text{ is a good approximation to } x) \geq 1 - \delta.$$



We have already seen nonuniform schemes in this course when discussing the point query problem in [Section 4.1](#). Recall in point query, when using a linear sketch  $\Pi$  we would like that given  $\Pi x$ , for any  $i$  we can recover an approximate  $\tilde{x}_i$  that is close to  $x_i$ ; if we set the failure probability to be  $\ll 1/n$ , then this approximation holds for all  $i \in [n]$  simultaneously with good probability. The **CountMin** sketch provides the  $\ell_\infty/\ell_1$  guarantee  $\|x - \tilde{x}\|_\infty \leq (1/k) \cdot \|x_{tail(k)}\|_\infty$ . The **CountSketch** provides the  $\ell_\infty, \ell_2$  guarantee  $\|x - \tilde{x}\|_\infty \leq (1/\sqrt{k}) \cdot \|x_{tail(k)}\|_2$ . These are both nonuniform schemes.

In this chapter, we will show how to obtain *uniform* schemes via measurements matrices that satisfy the restricted isometry property ([Definition 5.3.5](#)). For a more thorough introduction to compressed sensing, see the textbook by Foucart and Rauhut [[FR13](#)].

## 7.1 Basis Pursuit

If  $x$  is not just approximately sparse but actually  $k$ -sparse, we have already seen a deterministic (and thus uniform) scheme to recover  $x$  exactly: the  $k$ -sparse recovery data structure of [Subsection 4.2.1](#), which is actually a linear sketch. For this we required that the measurement matrix  $\Pi \in \mathbb{R}^{m \times n}$  satisfied the property that every  $m \times 2k$  submatrix was of full column rank, which is equivalent to  $\Pi_S^\top \Pi_S$  having no zero eigenvalues for any  $S \subset [n]$  of size  $2k$  (where  $\Pi_S$  is the  $|S| \times n$  matrix obtained by restricting  $\Pi$  to only contain the columns in  $S$ ). With such a matrix in hand, if  $x$  is  $k$ -sparse we can recover it exactly given  $\Pi x$  by solving the following optimization problem:

$$\begin{array}{ll} \min_{z \in \mathbb{R}^n} & \|z\|_0 \\ \text{s.t.} & \Pi z = \Pi x \end{array}$$

Unfortunately, this optimization problem is NP-hard to solve in general [[GJ79](#), Problem MP5], though for particular  $\Pi$  can be tractable (e.g. the polynomial-time “syndrome decoding” algorithm in the case that  $\Pi$  is as in [Subsection 4.2.1](#)).

What about the case though when  $x$  is not exactly  $k$ -sparse, but only approximately so? The  $(\varepsilon, k)$ -restricted isometry property (RIP) provides a robustification of this property: for any  $S$  of size  $k$  we not only require that  $\Pi_S^\top \Pi_S$  has no zero eigenvalues, but furthermore that all its eigenvalues lie in the interval  $[1 - \varepsilon, 1 + \varepsilon]$ . It turns out that if  $\Pi$  satisfies the RIP, there is a polynomial-time algorithm to find an  $\tilde{x}$  that well-approximates  $x$ , by solving a linear program [[CRT06](#), [Don06](#)]. This linear program is known as *basis pursuit*, and is the following:

$$\begin{array}{ll} \min_{z \in \mathbb{R}^n} & \|z\|_1 \\ \text{s.t.} & \Pi z = \Pi x \end{array}$$

Note this is a linear program since we can define variables  $y_1, \dots, y_n$  and add the constraints  $y_i \geq z_i, y_i \geq -z_i$  for each  $i$  then minimize  $\sum_i y_i$ .

We now show that if  $\Pi$  satisfies RIP, basis pursuit returns a good approximation to  $x$ .

**Theorem 7.1.1** ([[Can08](#)]). *If  $\Pi$  is  $(\varepsilon_{2k}, 2k)$ -RIP with  $\varepsilon_{2k} \leq \sqrt{2} - 1$ , and  $\tilde{x} = x + h$  is the optimal solution to the basis pursuit linear program, then*

$$\|h\|_2 \leq O\left(\frac{1}{\sqrt{k}}\right) \|x_{tail(k)}\|_1.$$

*Proof.* First, we define some notation.

For a vector  $x \in \mathbb{R}^n$  and set  $S \subseteq [n]$ , let  $x_S$  be the vector with all of its coordinates outside of  $S$  zeroed out. We use  $\bar{T}_i$  to indicate the complement of  $T_i$ . Let  $T_0 \subseteq [n]$  be the indices of the

largest (in absolute value)  $k$  coordinates of  $x$ . Then for  $j \geq 1$ , let  $T_j$  be the indices of the largest  $k$  coordinates (in absolute value) of  $h_{\cup_{i < j} T_i}$ .

By the triangle inequality,

$$\begin{aligned} \|h\|_2 &= \|h_{T_0 \cup T_1} + h_{\overline{T_0 \cup T_1}}\|_2 \\ &\leq \|h_{T_0 \cup T_1}\|_2 + \|h_{\overline{T_0 \cup T_1}}\|_2. \end{aligned}$$

Our strategy for bounding  $h$  will be to show:

- $\|h_{\overline{T_0 \cup T_1}}\|_2 \leq \|h_{T_0 \cup T_1}\|_2 + O\left(\frac{1}{\sqrt{k}}\right) \|x_{tail(k)}\|_1$ , and
- $\|h_{T_0 \cup T_1}\|_2 \leq O\left(\frac{1}{\sqrt{k}}\right) \|x_{tail(k)}\|_1$ .

Both parts rely on the following lemma.

**Lemma 7.1.2.**

$$\sum_{j \geq 2} \|h_{T_j}\|_2 \leq \frac{2}{\sqrt{k}} \|x_{T_0^c}\|_1 + \|h_{T_0 \cup T_1}\|_2.$$

*Proof.* We apply the shelling trick (recall Eq. (5.10)).

$$\begin{aligned} \sum_{j \geq 2} \|h_{T_j}\|_2 &\leq \sqrt{k} \sum_{j \geq 2} \|h_{T_j}\|_\infty && (h_{T_j} \text{ is } k\text{-sparse}) \\ &\leq \frac{1}{\sqrt{k}} \sum_{j \geq 2} \|h_{T_{j-1}}\|_1 && (\text{shelling}) \\ &\leq \frac{1}{\sqrt{k}} \|h_{\overline{T_0}}\|_1. && (7.3) \end{aligned}$$

Now since  $\tilde{x} = x + h$  is the minimizer of the LP, we must have

$$\begin{aligned} \|x\|_1 &\geq \|x + h\|_1 \\ &= \|(x + h)_{T_0}\|_1 + \|(x + h)_{\overline{T_0}}\|_1 \\ &\geq \|x_{T_0}\|_1 - \|h_{T_0}\|_1 + \|h_{\overline{T_0}}\|_1 - \|x_{\overline{T_0}}\|_1 \end{aligned}$$

by two applications of the triangle inequality. Rearranging,

$$\begin{aligned} \|h_{\overline{T_0}}\|_1 &\leq \|x\|_1 - \|x_{T_0}\|_1 + \|h_{T_0}\|_1 + \|x_{\overline{T_0}}\|_1 \\ &= 2\|x_{\overline{T_0}}\|_1 + \|h_{T_0}\|_1 \\ &\leq 2\|x_{\overline{T_0}}\|_1 + \sqrt{k} \|h_{T_0}\|_2 && (\text{Cauchy-Schwarz}) \\ &\leq 2\|x_{\overline{T_0}}\|_1 + \sqrt{k} \|h_{T_0 \cup T_1}\|_2 \end{aligned}$$

Combining this upper bound with Eq. (7.3) yields the claim.  $\square$

Returning to the main proof, let us first upper bound the size of  $h_{\overline{T_0 \cup T_1}}$ . We get:

$$\begin{aligned}
\|h_{\overline{T_0 \cup T_1}}\|_2 &= \left\| \sum_{j \geq 2} h_{T_j} \right\|_2 \\
&\leq \sum_{j \geq 2} \|h_{T_j}\|_2 \\
&\leq \|h_{T_0 \cup T_1}\|_2 + \frac{2}{\sqrt{k}} \|x_{\overline{T_0}}\|_1 \quad (\text{Lemma 7.1.2}) \\
&= \|h_{T_0 \cup T_1}\|_2 + \frac{2}{\sqrt{k}} \|x_{\text{tail}(k)}\|_1.
\end{aligned}$$

To bound the size of  $h_{T_0 \cup T_1}$ , first observe that

$$\Pi h_{T_0 \cup T_1} = \Pi h - \sum_{j \geq 2} \Pi h_{T_j} = - \sum_{j \geq 2} \Pi h_{T_j}$$

since  $h \in \ker \Pi$ . Therefore,

$$\|\Pi h_{T_0 \cup T_1}\|_2^2 = - \sum_{j \geq 2} \langle \Pi h_{T_0 \cup T_1}, \Pi h_{T_j} \rangle \leq \sum_{j \geq 2} (|\langle \Pi h_{T_0}, \Pi h_{T_j} \rangle| + |\langle \Pi h_{T_1}, \Pi h_{T_j} \rangle|).$$

By part (b) of [Lemma 5.3.6](#), each summand is at most

$$\begin{aligned}
\varepsilon_{2k} (\|h_{T_0}\|_2 + \|h_{T_1}\|_2) \|h_{T_j}\|_2 &= \varepsilon_{2k} \sqrt{(\|h_{T_0}\|_2 + \|h_{T_1}\|_2)^2} \|h_{T_j}\|_2 \\
&= \varepsilon_{2k} \sqrt{\|h_{T_0}\|_2^2 + \|h_{T_1}\|_2^2 + 2\|h_{T_0}\|_2 \|h_{T_1}\|_2} \|h_{T_j}\|_2 \\
&\leq \varepsilon_{2k} \sqrt{2\|h_{T_0}\|_2^2 + 2\|h_{T_1}\|_2^2} \|h_{T_j}\|_2 \quad (\text{AM-GM}) \\
&= \varepsilon_{2k} \sqrt{2} \|h_{T_0 \cup T_1}\|_2 \|h_{T_j}\|_2.
\end{aligned}$$

Thus

$$\begin{aligned}
(1 - \varepsilon_{2k}) \|h_{T_0 \cup T_1}\|_2^2 &\leq \|\Pi h_{T_0 \cup T_1}\|_2^2 \\
&\leq \varepsilon_{2k} \sqrt{2} \|h_{T_0 \cup T_1}\|_2 \sum_{j \geq 2} \|h_{T_j}\|_2 \\
&\leq \varepsilon_{2k} \sqrt{2} \|h_{T_0 \cup T_1}\|_2 \left( \frac{2}{\sqrt{k}} \|x_{\overline{T_0}}\|_1 + \|h_{T_0 \cup T_1}\|_2 \right) \quad (\text{Lemma 7.1.2})
\end{aligned}$$

Cancelling a factor of  $\|h_{T_0 \cup T_1}\|_2$  from both sides and rearranging gives

$$\|h_{T_0 \cup T_1}\|_2 \leq \frac{\varepsilon_{2k} 2 \sqrt{2}}{(1 - \varepsilon_{2k} - \varepsilon_{2k} \sqrt{2}) \sqrt{k}} \|x_{\overline{T_0}}\|_1 = O\left(\frac{1}{\sqrt{k}}\right) \|x_{\text{tail}(k)}\|_1,$$

with the last equality holding since  $\varepsilon_{2k} < \sqrt{2} - 1$ . Putting everything together:

$$\|h\|_2 \leq \|h_{\overline{T_0 \cup T_1}}\|_2 + \|h_{T_0 \cup T_1}\|_2$$

$$\begin{aligned} &\leq 2\|h_{T_0 \cup T_1}\|_2 + \frac{2}{\sqrt{k}}\|x_{tail(k)}\|_1 \\ &\leq O\left(\frac{1}{\sqrt{k}}\right)\|x_{tail(k)}\|_1. \end{aligned}$$

□

### 7.1.1 Obtaining RIP matrices

Thus far we have shown how to use RIP matrices in [Subsection 5.3.3](#) to obtain Fast JL, and in [Section 7.1](#) for compressed sensing. But how does one obtain RIP matrices? There are generally three ways:

- **Via the JL lemma/OSE's.** Note that being an RIP matrix means being a subspace embedding for  $\binom{n}{k}$  subspaces simultaneously (namely the set of all  $k$ -dimensional subspaces obtained spanned by choosing  $k$  distinct standard basis vectors). A matrix with i.i.d. subgaussian entries is an  $\varepsilon$ -OSE for  $d$ -dimensional subspaces with probability at least  $1 - \delta$  as long as its number of rows is at least  $\Omega(\varepsilon^{-2}(d + \log(1/\delta)))$ . Setting  $d = k$  and  $\delta < 1/\binom{n}{k}$  to union bound gives that  $m = O(\varepsilon^{-2}k \log(n/k))$  suffices.
- **Analyzing suprema of random processes.** There are other random matrices that people have studied in the context of RIP, such as sampling rows from the Hadamard matrix or Discrete Fourier Transform [[CT06](#), [RV08](#), [Bou14](#), [HR16](#)]. One then wishes to analyze  $\mathbb{E}_\Pi \sup_{x \in D_{n,k}} |||\Pi x\|_2^2 - 1|$  and show that it is at most  $\varepsilon$ , where  $D_{n,k}$  is the set of all unit norm vectors in  $\mathbb{R}^n$  which are  $k$ -sparse. Some of these arguments use analyses such as Dudley's inequality, generic chaining, and other tools from the study of the suprema of random processes.
- **Ad hoc.** Some analyses of RIP are “one-offs” in that the analysis was invented for the specific matrix being studied and not used again in the context of analyzing RIP. For example, motivated by an explicit, deterministic construction of RIP matrices, Bourgain et al. analyzed a construction based on tools from analytic number theory [[BDF<sup>+</sup>11](#)]. The construction achieves  $O(k^{2-\gamma}(\varepsilon^{-1} \log n)^c)$  rows for some absolute constant  $c > 0$  for a narrow range of  $k$  close to  $\sqrt{n}$ . This is the only known explicit construction of an RIP matrix to achieve a subquadratic number of rows in  $k$ .
- **Incoherent matrices.** A matrix  $\Pi \in \mathbb{R}^{m \times n}$  is said to be  $\alpha$ -incoherent if (1) its columns  $\pi_i$  each have unit Euclidean norm, and (2) for all  $i \neq j$ ,  $|\langle \pi_i, \pi_j \rangle| \leq \alpha$ . We show below that any such matrix with  $\alpha \leq \varepsilon/(k-1)$  is  $\alpha$ -incoherent. One can construct such matrices using Reed-Solomon codes with  $m = O(\alpha^{-2}((\log n)/(\log \log n + \log(1/\alpha)))^2)$ , or using random codes with  $m = O(\alpha^{-2} \log n)$ . Either of these leads to RIP matrices with  $\varepsilon^{-2}k^2 \log^c n$  rows.

**Theorem 7.1.3** (Gershgorin circle theorem). *Let  $\lambda_1, \dots, \lambda_N \in \mathbb{C}$  be the eigenvalues of a matrix  $A \in \mathbb{R}^{N \times N}$ . Then for each  $i \in [N]$ , there exists a  $j \in [N]$  such that  $\lambda_i$  lives in a complex disc about  $A_{j,j}$  of radius  $\sum_{r \neq j} |A_{j,r}|$ .*

*Proof.* Let  $v_i$  be the eigenvector corresponding to  $\lambda_i$ . Let  $j$  be such that  $\|v_i\|_\infty = |(v_i)_j|$ . Then  $(Av_i)_j = \lambda_i(v_i)_j$ , but also

$$(Av_i)_j = A_{j,j}(v_i)_j + \sum_{r \neq j} A_{r,j}(v_i)_r.$$

Thus

$$(v_i)_j(A_{j,j} - \lambda_i) = \sum_{r \neq j} A_{r,j}(v_i)_r,$$

implying

$$|A_{j,j} - \lambda_i| = \left| \sum_{r \neq j} A_{r,j} \frac{(v_i)_r}{(v_i)_j} \right| \leq \sum_{r \neq j} \left| A_{r,j} \frac{(v_i)_r}{(v_i)_j} \right| \leq \sum_{r \neq j} |A_{r,j}|.$$

□

**Corollary 7.1.4.** *If  $\Pi \in \mathbb{R}^{m \times n}$  is  $\alpha$ -incoherent for  $\alpha \leq \varepsilon/(k-1)$ , then  $\Pi$  satisfies  $(\varepsilon, k)$ -RIP.*

*Proof.* Let  $S \subset [n]$  be of size  $k$ . We would like that all the eigenvalues of  $A_S := \Pi_S^\top \Pi_S$  lie in the interval  $[1 - \varepsilon, 1 + \varepsilon]$ , where  $\Pi_S \in \mathbb{R}^{m \times k}$  is the restriction of  $\Pi$  to columns in  $S$ . Note all diagonal entries of  $A_S$  are 1 and the off-diagonals are at most  $\alpha$  in magnitude due to incoherence. Furthermore all the eigenvalues are real by the spectral theorem, since  $A_S$  is a real, symmetric matrix. Thus by the Gershgorin circle theorem, all eigenvalues of  $A_S$  are in the interval  $[1 - \alpha(k-1), 1 + \alpha(k-1)]$ , as desired. □

## 7.2 Iterative Hard Thresholding

Though Basis Pursuit provides a uniform recovery guarantee, we would prefer to not use a generic LP solver as that would be slow. Thankfully, there are fast *iterative* approaches developed for (approximately) recovering  $x$  from measurements which have running time nearly linear (if the measurement matrix supports nearly linear time matrix-vector multiplication, such as sampling rows from the DFT). An iterative approach called CoSAMP was first introduced by Needell and Tropp [NT08]. In this section we cover a different algorithm, *Iterative Hard Thresholding (IHT)* and it is due to Blumensath and Davies [BD09].

---

**Algorithm 1** Iterative Hard Thresholding (IHT).

---

```

1: function IHT( $\Pi, y(= \Pi x + e), k, T$ )
2:    $x^{[0]} \leftarrow 0$ 
3:   for  $t = 0 \dots T - 1$  do
4:      $x^{[t+1]} \leftarrow H_k(x^{[t]} + \Pi^\top (y - \Pi x^{[t]})) \triangleright H_k(z) := z_{head(k)}$ 
5:   end for
6:   return  $x^{[T]}$ 
7: end function

```

---

The IHT algorithm starts with initial iterate  $x^{[0]} = 0$  then iteratively improves it to move closer to  $x$ . Suppose the iterates produced over  $T$  iterations are  $x^{[1]}, \dots, x^{[T]}$ . Then the main theorem of IHT is as follows.

**Theorem 7.2.1** ([BD09]). *If  $\Pi$  satisfies  $(\varepsilon, 3k)$ -RIP for  $\varepsilon < \frac{1}{4\sqrt{2}}$ , then  $\forall T \geq 1$*

$$\|x^{[T]} - x\|_2 \lesssim 2^{-T} \|x\|_2 + \|x_{tail(k)}\|_2 + \frac{1}{\sqrt{k}} \|x_{tail(k)}\|_1 + \|e\|_2 \quad (7.4)$$

Here we consider a more general problem than Section 7.1; specifically we assume there may be some *post-measurement noise*  $e$ , so that what we obtain from our measurements is actually not

$\Pi x$ , but rather  $\Pi x + e$ . For example, the measurements may be made by a physical device, and our measurement readings will then have error introduced by finite precision and other potential sources of noise.

Comparing to the guarantee of basis pursuit, Eq. (7.4) we have three extra terms:  $2^{-T}\|x\|_2$ ,  $\|x_{tail(k)}\|_2$ , and  $\|e\|_2$ . Note that the last term corresponds to the post-measurement noise and it is unavoidable. For the second term,  $\|x_{tail(k)}\|_2$ , we show (Claim 7.2.2) that it is dominated by  $O(\|x_{tail(k/2)}\|_1/\sqrt{k})$ . Hence, the only difference is the exponentially decaying term  $2^{-T}\|x\|_2$ . In turn, the IHT algorithm is much faster than using an off-the-shelf LP solver to solve basis pursuit.

**Claim 7.2.2.**  $\|x_{tail(2k)}\|_2 \leq \frac{1}{\sqrt{k}}\|x_{tail(k)}\|_1$ .

*Proof.* Without loss of generality assume  $|x_1| \geq |x_2| \geq \dots \geq |x_n|$ . We partition the coordinates of  $x$  into blocks of size  $k$  with  $B_j = \{(n/k)j + 1, \dots, (n/k)j + k\}$ . Now we apply shelling (recall Eq. (5.10)):

$$\begin{aligned} \|x_{tail(2k)}\|_2^2 &= \sum_{j \geq 3} \|x_{B_j}\|_2^2 \\ &\leq \sum_{j \geq 3} k \cdot \|x_{B_j}\|_\infty^2 \\ &\leq \sum_{j \geq 3} k \cdot \left( \frac{\|x_{B_{j-1}}\|_1}{k} \right)^2 && \text{(shelling)} \\ &= \sum_{j \geq 2} \frac{1}{k} \cdot \|x_{B_j}\|_1^2 \end{aligned}$$

Now take square roots and use that  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ , noting that  $\sum_{j \geq 2} \|x_{B_j}\|_1 = \|x_{tail(k)}\|_1$ .  $\square$

Now we focus on the proof of the convergence of IHT (proof of Theorem 7.2.1). Note in the analysis we can assume that  $x$  is *exactly*  $k$ -sparse. More precisely, we decompose  $x = x_{head(k)} + x_{tail(k)}$  so that  $\Pi x + e = \Pi x_{head(k)} + (\Pi x_{tail(k)} + e) = \Pi x_{head(k)} + \tilde{e}$ , defining  $\tilde{e}$  as  $\Pi x_{tail(k)} + e$ . Also, subdivide each  $B_j$  into  $C_j, C'_j$  each of size  $k/2$ , where  $C_j$  contains the first  $k/2$  elements of  $B_j$  and  $C'_j$  contains the next  $k/2$  elements. Then

$$\begin{aligned} \|\tilde{e}\|_2 &\leq \|e\|_2 + \|\Pi x_{tail(k)}\|_2 \\ &= \|e\|_2 + \left\| \sum_{j \geq 2} \Pi x_{B_j} \right\|_2 \\ &\leq \|e\|_2 + \sum_{j \geq 2} \|\Pi x_{B_j}\|_2 \\ &\leq \|e\|_2 + (1 + \varepsilon) \sum_{j \geq 2} \|x_{B_j}\|_2 && \text{(RIP)} \\ &\leq \|e\|_2 + (1 + \varepsilon) \sum_{j \geq 2} \sqrt{k} \cdot \|x_{B_j}\|_\infty \\ &\leq \|e\|_2 + (1 + \varepsilon) \sum_{j \geq 2} \sqrt{k} \cdot \left( \frac{\|x_{C'_{j-1}}\|_1}{k/2} \right) \\ &\leq \|e\|_2 + \frac{2(1 + \varepsilon)}{\sqrt{k}} \|x_{tail(k/2)}\|_1, \end{aligned}$$

which is an allowable error for [Theorem 7.2.1](#).

We now prove [Theorem 7.2.1](#).

*Proof.* We measure the progress of IHT based on the residual vector  $r^{[t]} := x - x^{[t]}$ . We show that  $\|r^{[t]}\|_2$  decreases at some rate as  $t$  increases. For analysis purposes, we define  $a^{[t+1]} := x^{[t]} + \Pi^\top(y - \Pi x^{[t]})$  (note that  $x^{[t+1]} = H_k(a^{[t+1]})$ ).

We make the following definitions:

- $\Gamma_k^* = \text{supp}(x)$ ,
- $\Gamma^{[t]} = \text{supp}(x^{[t]})$ , and
- $B^{[t]} = \Gamma_k^* \cup \Gamma^{[t]}$ .

We now bound  $\|r^{[t+1]}\|_2$ . Below let  $B$  denote  $B^{[t+1]}$  and  $B' = B^{[t]}$ .

$$\begin{aligned}
\|r^{[t+1]}\|_2 &= \|x - x^{[t+1]}\|_2 \\
&= \|x_B - x_B^{[t+1]}\|_2 \\
&\leq \|x_B - a_B^{[t+1]}\|_2 + \|a_B^{[t+1]} - x_B^{[t+1]}\|_2 \\
&\leq 2\|x_B - a_B^{[t+1]}\|_2 && (x^{[t+1]} \text{ is the best } k\text{-sparse approx. to } a^{[t+1]}) \\
&= 2\|x_B - x_B^{[t]} - \Pi_B^\top(y - \Pi x^{[t]})\|_2 \\
&= 2\|r_B^{[t]} - \Pi_B^\top(\Pi r^{[t]} + e)\|_2 \\
&= 2\|r_B^{[t]} - \Pi_B^\top(\Pi_B r_B^{[t]} + \Pi_{B' \setminus B} r_{B' \setminus B}^{[t]} + e)\|_2 && (r^{[t]} = r_B^{[t]} + r_{B' \setminus B}^{[t]}) \\
&\leq 2\|(I_B - \Pi_B^\top \Pi_B) r_B^{[t]}\|_2 + 2\|\Pi_B^\top \Pi_{B' \setminus B} r_{B' \setminus B}^{[t]}\|_2 + 2\|\Pi_B^\top e\|_2 && (\text{triangle inequality}) \\
&\leq 2\|I_B - \Pi_B^\top \Pi_B\| \cdot \|r_B^{[t]}\|_2 + 2\|\Pi_B^\top \Pi_{B' \setminus B}\| \cdot \|r_{B' \setminus B}^{[t]}\|_2 + 2\|\Pi_B^\top\| \cdot \|e\|_2 \\
&\leq 2\varepsilon(\|r_B^{[t]}\|_2 + \|r_{B' \setminus B}^{[t]}\|_2) + 2\sqrt{1 + \varepsilon}\|e\|_2 && (\text{Lemma 5.3.6}) \\
&\leq 2\sqrt{2}\varepsilon\|r^{[t]}\|_2 + 3\|e\|_2 && (7.5) \\
&\leq \frac{1}{2}\|r^{[t]}\|_2 + 3\|e\|_2, && (7.6)
\end{aligned}$$

[Eq. \(7.6\)](#) holds since  $\varepsilon \leq 1/(4\sqrt{2})$ . [Eq. \(7.5\)](#) follows from AM-GM. To see this, write  $\alpha = \|r_B^{[t]}\|_2$ ,  $\beta = \|r_{B' \setminus B}^{[t]}\|_2$ ,  $\gamma = \|r^{[t]}\|_2$ , and note  $\gamma = \sqrt{\alpha^2 + \beta^2}$ . Then

$$\begin{aligned}
(\alpha + \beta)^2 &= (\alpha^2 + \beta^2) + 2\alpha\beta \\
&\leq (\alpha^2 + \beta^2) + (\alpha^2 + \beta^2) && (\text{AM-GM}) \\
&= 2\gamma^2,
\end{aligned}$$

so that  $\alpha + \beta \leq \sqrt{2}\gamma$ , as claimed.

We then have from [Eq. \(7.6\)](#) that  $\|r^{[t+1]}\|_2 \leq (1/2)\|r^{[t]}\|_2 + 3\|e\|_2$ . Since  $\|r^{[0]}\|_2 = \|x\|_2$ , an induction on  $t$  shows that for  $t \geq 1$ ,  $\|r^{[t]}\|_2 \leq 2^{-t}\|x\|_2 + 3(\sum_{i=0}^{t-1} 2^{-i})\|e\|_2$ , which is always at most  $2^{-t}\|x\|_2 + 6\|e\|_2$ .





## Chapter 8

# Suprema of stochastic processes and applications

A stochastic process is simply a collection of random variables,  $\{X_t\}_{t \in T}$ . Often one considers such processes that either have some temporal structure (i.e. the process evolves over time, such as Brownian motion), or spatial structure so that there is some important metric structure on the random variables that explains correlations between them, e.g. if one defines a distance such as  $d(X, Y) := (\mathbb{E}(X - Y)^2)^{1/2}$ . In this chapter we focus on understanding  $\mathbb{E} \sup_{t \in T} X_t$ , and using this understanding for sketching applications. We focus on spatial processes, where typically the underlying metric structure helps us control the expected supremum.

Because of the applications we consider, we will be focused on (sub)gaussian processes. Here there is some  $T \subset \mathbb{R}^d$ , and  $X_t$  for  $t \in T$  is simply  $\langle g, t \rangle$  for a gaussian vector  $g$  with mean zero and identity covariance matrix. We have the following definition:

**Definition 8.0.1.** For a subset  $T \subset \mathbb{R}^d$ , we define the *gaussian mean width*  $w(T)$  by

$$w(T) := \mathbb{E} \sup_{g, x \in S} \langle g, x \rangle,$$

where  $g$  is a  $d$ -dimensional gaussian with mean 0 and identity covariance matrix.

### 8.1 Methods of bounding gaussian mean width

We present four different ways of bounding the mean width  $w(T)$  of a set  $T \subset B_{\ell_2^d}$ , the unit Euclidean ball in  $\mathbb{R}^d$ . The methods presented are increasingly sophisticated.

**Union bound.** The first observation is that for  $x \in T$ ,  $\langle g, x \rangle$  is a gaussian with variance  $\|x\|_2^2 \leq 1$ . Thus  $(\langle g, x \rangle)_{x \in T}$  is a collection of  $|T|$  gaussians with variance at most 1, and thus we expect the maximum to be  $O(\sqrt{\log T})$ . Details follow.

$$\begin{aligned} \mathbb{E} \sup_{g, x \in T} \langle g, x \rangle &\leq \int_0^\infty \mathbb{P}(\sup_{x \in T} \langle g, x \rangle > t) dt \\ &= \underbrace{\int_0^{t^*} \mathbb{P}(\sup_{x \in T} \langle g, x \rangle > t) dt}_{\leq 1} + \int_{t^*}^\infty \mathbb{P}(\sup_{x \in T} \langle g, x \rangle > t) dt \end{aligned}$$

$$\begin{aligned} &\leq t^* + \sum_{x \in T} \int_{t^*}^{\infty} \mathbb{P}(|\langle g, x \rangle| > t) dt \\ &\leq t^* + 2|T| \cdot e^{-t^2/2}, \end{aligned}$$

which is  $O(\sqrt{\log |T|})$  by choosing  $t^* = c\sqrt{\log |T|}$ .

**$\varepsilon$ -net.** Let  $T' \subseteq T$  be an  $\varepsilon$ -net of  $T$  under  $\ell_2$ . Then for any  $x \in T$  let  $x' \in T'$  be the closest element in  $T'$  to  $x$ . Then

$$\begin{aligned} \mathbb{E} \sup_{x \in T} \langle g, x \rangle &\leq \mathbb{E} \sup_{x \in T'} \langle g, x' \rangle + \mathbb{E} \sup_{x \in T} \langle g, x - x' \rangle \\ &\lesssim \sqrt{\log |T'|} + \left( \sup_{x \in T} \|x - x'\|_2 \right) \cdot \mathbb{E} \sup_{x \in T'} \left\langle g, \frac{x - x'}{\|x - x'\|_2} \right\rangle \\ &\lesssim \sqrt{\log |T'|} + \varepsilon \sqrt{\log |T|} \end{aligned}$$

Since  $\varepsilon$  can be picked arbitrarily, we have the bound  $w(T) \lesssim \inf_{\varepsilon \geq 0} \{ \varepsilon \sqrt{\log |T|} + \log^{1/2} \mathcal{N}(T, \ell_2, \varepsilon) \}$ . Note this is more general than the union bound argument, which we simply recover by picking  $\varepsilon = 0$ . We also note that if  $|T|$  is infinite, we can alternatively bound  $\mathbb{E} \sup_{x \in T} \langle g, x - x' \rangle$  by  $\varepsilon \sqrt{d}$  using Cauchy-Schwarz.

**Dudley's inequality.** In Dudley's inequality we pick not one  $\varepsilon$ -net, but an infinite sequence of nets  $(S_k)_{k=0}^{\infty}$  where  $\varepsilon_k = 2^{-k}$ . Without loss of generality we can assume  $|T|$  is finite, since otherwise we can first approximate  $T'$  by an arbitrarily fine net (arbitrarily small  $\varepsilon$ ) as discussed above, then apply Dudley's inequality to  $T'$ . Let  $\pi_k x$  denote the closest point in  $S_k$  to  $x$ , and define  $\Delta_k x := \pi_k x - \pi_{k-1} x$ . Note  $S_0$  can be taken as  $\{0\}$ , since  $T$  is a subset of the unit ball. Then  $x = \sum_{k=1}^{\infty} \Delta_k x$  so that

$$\begin{aligned} \mathbb{E} \sup_{x \in T} \langle g, x \rangle &\leq \sum_{k=1}^{\infty} \mathbb{E} \sup_{x \in T} \langle g, \Delta_k x \rangle \\ &\leq \sum_{k=1}^{\infty} \left( \sup_{x \in T} \|\Delta_k x\|_2 \right) \cdot \mathbb{E} \sup_{x \in T} \left\langle g, \frac{\Delta_k x}{\|\Delta_k x\|_2} \right\rangle \\ &\lesssim \sum_{k=1}^{\infty} \left( \sup_{x \in T} (\|\pi_k x\|_2 + \|\pi_{k-1} x\|_2) \right) \cdot \log^{1/2} \left( \mathcal{N}(T, \ell_2, 1/2^k) \cdot \mathcal{N}(T, \ell_2, 1/2^{k-1}) \right) \\ &\lesssim \sum_{k=1}^{\infty} \frac{1}{2^k} \cdot \log^{1/2} \mathcal{N}(T, \ell_2, \frac{1}{2^k}) \end{aligned} \tag{8.1}$$

The final bound above on  $w(T)$  is known as *Dudley's inequality*, and since the covering number of the unit ball of an  $r$ -dimensional subspace of  $\mathbb{R}^d$  under  $\ell_2$  is at most  $\sqrt{r}$  it implies together with Gordon's theorem that a bound of  $m = O(r/\varepsilon^2)$  suffices for oblivious subspace embeddings (see [Subsection 6.2.3](#)).

Note: the right hand side of Dudley's inequality is often in the literature bounded by an integral. That is, many authors write Dudley's inequality as

$$w(T) \lesssim \int_0^{\infty} \log^{1/2} \mathcal{N}(T, \ell_2, u) du.$$

**Generic chaining.** Let  $T$  be a finite subset of some normed vector space with norm  $\|\cdot\|_X$ . We say that a sequence  $T_0 \subseteq T_1 \subseteq \dots \subseteq T$  is *admissible* if  $|T_0| = 1$  and  $|T_r| \leq 2^{2^r}$  for all  $r \geq 1$ , and  $T_r = T$  for all  $r \geq r_0$  for some  $r_0$ . We define the  $\gamma_2$ -functional

$$\gamma_2(T, \|\cdot\|_X) = \inf \sup_{x \in T} \sum_{r=1}^{\infty} 2^{r/2} \cdot d_X(x, T_r),$$

where the inf is taken over all admissible sequences.

The first observation is that the bound from Dudley's inequality can be rewritten in an equivalent form that only differs by a constant factor:

$$w(T) \lesssim \inf_{\{T_r\}} \sum_{r=0}^{\infty} 2^{r/2} \cdot \sup_{x \in T} d_{\ell_2}(x, T_r),$$

where  $\{T_r\}$  is admissible. To see why, first observe that to optimize the choice of the  $T_r$ , we should also do something like compute  $\varepsilon$ -nets: given a budget of  $2^{2^r}$  for the size of  $T_r$ , we should choose  $T_r$  so that it is an  $\varepsilon$ -net for the smallest value of  $\varepsilon$  possible under this size constraint. Note then that the above supremum is  $\in (1/2, 1]$  for  $r = 0, 1, \dots, r_1 - 1$  for some  $r_1$ , and in general is in  $(2^{-(k+1)}, 2^{-k}]$  for  $r = r_k, \dots, r_{k+1} - 1$ . During each of these "phases" (different  $k$ ), the suprema above over  $x \in T$  are the same up to a factor of 2, while the summation over  $r$  of  $2^{r/2}$  is a geometric series dominated by its last term and is thus  $O(2^{r_{k+1}/2})$ . Thus  $2^{-k}$  is essentially as in Eq. (8.1) with the  $O(2^{r_{k+1}/2})$  term serving the role of  $\log^{1/2} \mathcal{N}(T, \ell_2, 1/2^k)$ .

A surprising result of Fernique, developed further by Talagrand, is that the inequality still holds even if the supremum over  $x$  is taken *outside* the summation, i.e.

$$w(T) \lesssim \inf_{\{T_r\}} \sup_{x \in T} \sum_{r=0}^{\infty} 2^{r/2} \cdot d_{\ell_2}(x, T_r).$$

The right hand side of the above is known as  $\gamma_2(T, \ell_2)$ . We know prove that this is the case.

Henceforth, as in the proof of Dudley's inequality we let  $\pi_r x$  denote the closest point to  $x$  in  $T_r$ , and  $\Delta_r x := \pi_r x - \pi_{r-1} x$ .

**Theorem 8.1.1.** *For any  $T$ ,  $w(T) \lesssim \gamma_2(T, \ell_2)$ .*

*Proof.* For  $x \in T$  we can write  $x = \pi_0 x + \sum_{r=1}^{\infty} \Delta_r x$ . Therefore

$$\begin{aligned} w(T) &\leq \mathbb{E} \sup_{x \in T} \langle g, \pi_0 x \rangle + \mathbb{E} \sup_{x \in T} \sum_{r=1}^{\infty} \langle g, \Delta_r x \rangle \\ &= \mathbb{E} \langle g, x_0 \rangle + \mathbb{E} \sup_{x \in T} \sum_{r=1}^{\infty} \langle g, \Delta_r x \rangle \quad (\text{since } |T_0| = 1, \text{ so } T_0 = \{x_0\} \text{ for some } x_0) \\ &= \mathbb{E} \sup_{x \in T} \sum_{r=1}^{\infty} \langle g, \Delta_r x \rangle \end{aligned}$$

Define  $X_r(x) := \langle g, \Delta_r x \rangle$ . Then by gaussian tail bounds, for any  $r$  and  $x$  we have

$$\mathbb{P}(|X_r(x)| > t 2^{r/2} \|\Delta_r x\|_2) < 2e^{-t^2 2^r}. \quad (8.2)$$

Thus

$$\begin{aligned}
\mathbb{E} \sup_{x \in T} \sum_{r=1}^{\infty} \langle g, \Delta_r x \rangle &\leq \int_0^{\infty} \mathbb{P} \left( \sup_{x \in T} \sum_{r=1}^{\infty} |X_r(x)| > u \right) du \\
&= \left( \sup_{x \in T} \sum_{r=1}^{\infty} 2^{r/2} \|\Delta_r x\|_2 \right) \cdot \int_0^{\infty} \mathbb{P} \left( \sup_{x \in T} \sum_{r=1}^{\infty} |X_r(x)| > t \cdot \sup_{x \in T} \sum_{r=1}^{\infty} 2^{r/2} \|\Delta_r x\|_2 \right) dt \\
&\leq \left( \sup_{x \in T} \sum_{r=1}^{\infty} 2^{r/2} \|\Delta_r x\|_2 \right) \cdot \left[ 3 + \int_3^{\infty} \mathbb{P} \left( \sup_{x \in T} \sum_{r=1}^{\infty} |X_r(x)| > t \cdot \sup_{x \in T} \sum_{r=1}^{\infty} 2^{r/2} \|\Delta_r x\|_2 \right) dt \right] \\
&\leq \left( \sup_{x \in T} \sum_{r=1}^{\infty} 2^{r/2} \|\Delta_r x\|_2 \right) \cdot \left[ 3 + \int_3^{\infty} \sum_{r=1}^{\infty} 2e^{-t^2/2} |T_r| |T_{r-1}| \right] \\
&\hspace{15em} \text{(Eq. (8.2) and union bound over all } r \text{ and } \Delta_r x) \\
&< \left( \sup_{x \in T} \sum_{r=1}^{\infty} 2^{r/2} \|\Delta_r x\|_2 \right) \cdot \left[ 3 + \int_3^{\infty} \sum_{r=1}^{\infty} 2e^{-t^2 2^r} (2^{2^r})^2 dt \right] \hspace{5em} \text{(since } |T_r| \leq 2^{2^r}) \\
&\lesssim \sup_{x \in T} \sum_{r=1}^{\infty} 2^{r/2} \|\Delta_r x\|_2 \\
&\leq \sup_{x \in T} \sum_{r=1}^{\infty} 2^{r/2} \cdot (d_{\ell_2}(x, T_r) + d_{\ell_2}(x, T_{r-1})) \hspace{5em} \text{(triangle inequality)} \\
&\lesssim \sup_{x \in T} \sum_{r=1}^{\infty} 2^{r/2} \cdot d_{\ell_2}(x, T_r)
\end{aligned}$$

Since we can choose  $\{T_r\}$  to be any admissible sequence, we may choose it to minimize the above expression, obtaining an upper bound on  $w(T)$  of  $\gamma_2(T, \ell_2)$ , as desired.  $\square$

## 8.2 Instance-wise bounds for Johnson-Lindenstrauss

Consider the following setup: we have some  $T \subseteq S^{d-1}$  (the unit Euclidean sphere in  $\mathbb{R}^d$ , i.e.  $S^{d-1} := \{x \in \mathbb{R}^n : \|x\|_2 = 1\}$ ), and we would like some matrix  $\Pi \in \mathbb{R}^{m \times d}$  to satisfy

$$\sup_{x \in T} \left| \|\Pi x\|_2^2 - 1 \right| \leq \varepsilon,$$

or if  $\Pi$  is random, that

$$\mathbb{E} \sup_{x \in T} \left| \|\Pi x\|_2^2 - 1 \right| \leq \varepsilon.$$

We can model the situation of the JL lemma, where we would like a low distortion embedding of  $X \subset \ell_2$  into  $\mathbb{R}^m$ ,  $|X| = n$ , by setting  $T = \{(x - y) / \|x - y\|_2 : x, y \in X\}$ . By choosing  $\Pi$  from an appropriate distribution, we know by the DJL lemma that  $m = O(\varepsilon^{-2} \log |T|)$  suffices. When  $T$  is  $S^{n-1}$  intersected with an  $r$ -dimensional linear subspace of  $\mathbb{R}^d$ ,  $m = O(r/\varepsilon^2)$  suffices ([Subsection 6.2.3](#)). Also, a  $\Pi$  with i.i.d. subgaussian entries satisfies  $(\varepsilon, k)$ -RIP for  $m = O(\varepsilon^{-2} k \log(ed/k))$ , corresponding to  $T$  being the set of all unit norm vectors with Euclidean norm exactly 1. What should  $m$  be in general, as a function of  $T$ ? Gordon gave an answer [[Gor88](#)] in terms of  $w(T)$ .

**Theorem 8.2.1** ([[Gor88](#)]). *Suppose  $\Pi \in \mathbb{R}^{m \times d}$  has  $\Pi_{i,j} \sim \mathcal{N}(0, 1/m)$  independent. There exists  $C > 0$  such that for any  $T \subseteq S^{d-1}$ ,  $m \geq C\varepsilon^{-2}(w^2(T) + 1)$  implies  $\mathbb{E}_{\Pi} \sup_{x \in T} \left| \|\Pi x\|_2^2 - 1 \right| \leq \varepsilon$ .*

The goal of this section is to prove (a version of) Gordon's theorem; specifically, assuming the *majorizing measures theorem* [Tal96], whose proof we do not give here, we will show that the statement of Gordon's theorem holds for  $\Pi_{i,j}$  being independent  $\pm 1/\sqrt{m}$  [KM05] (in fact any subgaussian distribution suffices, but we only provide the proof for Rademachers).

We prove [Theorem 8.2.1](#) using the following theorem of Kraher, Mendelson, and Rauhut [KMR14]. Below  $\rho_X(\mathcal{A})$  denotes the radius of  $\mathcal{A}$  under norm  $\|\cdot\|_X$ . We use  $\|\cdot\|_F$  to denote Frobenius norm, and  $\|\cdot\|$  to denote either  $\ell_2$  norm (for vectors) or  $\ell_2 \rightarrow \ell_2$  operator norm (for matrices).

**Theorem 8.2.2.** *Let  $\mathcal{A} \subset \mathbb{R}^{q \times q}$  be arbitrary and  $\sigma_1, \dots, \sigma_q$  be independent, uniform in  $\{-1, 1\}$ . Then*

$$\mathbb{E} \sup_{\sigma \in \mathcal{A}} \left| \|\mathcal{A}\sigma\|^2 - \mathbb{E} \|\mathcal{A}\sigma\|^2 \right| \lesssim \gamma_2^2(\mathcal{A}, \|\cdot\|) + \gamma_2(\mathcal{A}, \|\cdot\|) \cdot \rho_F(\mathcal{A}) + \rho_F(\mathcal{A}) \cdot \rho_{\ell_2 \rightarrow 2}(\mathcal{A}).$$

The KMR theorem was actually more general, where the Rademacher variables could be replaced by subgaussian random variables. We present just the proof of the Rademacher case.

*Proof.* Without loss of generality we can assume  $\mathcal{A}$  is finite (else apply the theorem to a sufficiently fine net, i.e. fine in  $\ell_2 \rightarrow \ell_2$  operator norm). Define

$$E = \mathbb{E} \sup_{\sigma \in \mathcal{A}} \left| \|\mathcal{A}\sigma\|^2 - \mathbb{E} \|\mathcal{A}\sigma\|^2 \right|$$

and let  $A^i$  denote the  $i$ th column of  $A$ . Then by decoupling ([Lemma 1.1.11](#))

$$\begin{aligned} E &= \mathbb{E} \sup_{\sigma \in \mathcal{A}} \left| \sum_{i \neq j} \sigma_i \sigma_j \langle A^i, A^j \rangle \right| \\ &\leq 4 \cdot \mathbb{E} \sup_{\sigma, \sigma' \in \mathcal{A}} \left| \sum_{i, j} \sigma_i \sigma'_j \langle A^i, A^j \rangle \right| \\ &= 4 \cdot \mathbb{E} \sup_{\sigma, \sigma' \in \mathcal{A}} |\langle \mathcal{A}\sigma, \mathcal{A}\sigma' \rangle|. \end{aligned}$$

Let  $\{T_r\}_{r=0}^\infty$  be admissible for  $\mathcal{A}$ . Direct computation shows

$$\langle \mathcal{A}\sigma, \mathcal{A}\sigma' \rangle = \langle (\pi_0 A)\sigma, (\pi_0 A)\sigma' \rangle + \sum_{r=1}^{\infty} \underbrace{\langle (\Delta_r A)\sigma, (\pi_{r-1} A)\sigma' \rangle}_{X_r(A)} + \sum_{r=1}^{\infty} \underbrace{\langle (\pi_r A)\sigma, (\Delta_r A)\sigma' \rangle}_{Y_r(A)}.$$

We have  $T_0 = \{A_0\}$  for some  $A_0 \in \mathcal{A}$ . Thus  $\mathbb{E}_{\sigma, \sigma'} |\langle (\pi_0 A)\sigma, (\pi_0 A)\sigma' \rangle|$  equals

$$\mathbb{E}_{\sigma, \sigma'} \left| \sigma^\top A_0^\top A_0 \sigma' \right| \leq \left( \mathbb{E}_{\sigma, \sigma'} \left( \sigma^\top A_0^\top A_0 \sigma' \right)^2 \right)^{1/2} = \|A_0^\top A_0\|_F \leq \|A_0\|_F \|A_0\| \leq \rho_F(\mathcal{A}) \cdot \rho_{\ell_2 \rightarrow 2}(\mathcal{A}).$$

Thus,

$$\mathbb{E}_{\sigma, \sigma'} \sup_{A \in \mathcal{A}} |\langle \mathcal{A}\sigma, \mathcal{A}\sigma' \rangle| \leq \rho_F(\mathcal{A}) \cdot \rho_{\ell_2 \rightarrow 2}(\mathcal{A}) + \mathbb{E}_{\sigma, \sigma'} \sup_{A \in \mathcal{A}} \sum_{r=1}^{\infty} |X_r(A)| + \mathbb{E}_{\sigma, \sigma'} \sup_{A \in \mathcal{A}} \sum_{r=1}^{\infty} |Y_r(A)|.$$

We focus on the first summation; handling the second summation over  $Y_r(A)$  is similar.

Note  $X_r(A) = \langle (\Delta_r A)\sigma, (\pi_{r-1}A)\sigma' \rangle = \langle \sigma, (\Delta_r A)^\top (\pi_{r-1}A)\sigma' \rangle$ . Thus

$$\mathbb{P}(|X_r(A)| > t2^{r/2} \cdot \|(\Delta_r A)^\top (\pi_{r-1}A)\sigma'\|) \leq 2e^{-t^22^{r/2}} \quad . \quad (\text{Khintchine})$$

Let  $\mathcal{E}(A)$  be the event that for all  $r \geq 1$  simultaneously,  $|X_r(A)| \leq t2^{r/2} \cdot \|\Delta_r A\| \cdot \sup_{A \in \mathcal{A}} \|A\sigma'\|$ . Then

$$\begin{aligned} \mathbb{P}(\exists A \in \mathcal{A} \text{ s.t. } \neg \mathcal{E}(A)) &\lesssim \sum_{r=1}^{\infty} |T_r| \cdot |T_{r-1}| \cdot e^{-t^22^{r/2}} \\ &\leq \sum_{r=1}^{\infty} 2^{2r+1} \cdot e^{-t^22^{r/2}}. \end{aligned}$$

Therefore

$$\mathbb{E}_{\sigma, \sigma'} \sup_{A \in \mathcal{A}} \sum_{r=1}^{\infty} |X_r(A)| = \mathbb{E}_{\sigma'} \int_0^{\infty} \mathbb{P}_{\sigma} \left( \sup_{A \in \mathcal{A}} \sum_{r=1}^{\infty} |X_r(A)| > t \right) dt,$$

which by a change of variables is equal to

$$\begin{aligned} &\mathbb{E}_{\sigma'} \left( \sup_{A \in \mathcal{A}} \|A\sigma'\| \cdot \left( \sup_{A \in \mathcal{A}} \sum_{r=1}^{\infty} 2^{r/2} \|\Delta_r A\| \right) \right. \\ &\quad \times \left. \int_0^{\infty} \mathbb{P}_{\sigma} \left( \sup_{A \in \mathcal{A}} \sum_{r=1}^{\infty} |X_r(A)| > t \sup_{A \in \mathcal{A}} \sum_{r=1}^{\infty} 2^{r/2} \cdot \|\Delta_r A\| \cdot \sup_{A \in \mathcal{A}} \|A\sigma'\| \right) dt \right) \\ &\leq \left( \mathbb{E}_{\sigma'} \sup_{A \in \mathcal{A}} \|A\sigma'\| \right) \cdot \left( \sup_{A \in \mathcal{A}} \sum_{r=1}^{\infty} 2^{r/2} \|\Delta_r A\| \right) \cdot \left[ 3 + \sum_{r=1}^{\infty} \int_3^{\infty} 2^{2r+1} e^{-t^22^{r/2}} dt \right] \\ &\lesssim \left( \mathbb{E}_{\sigma'} \sup_{A \in \mathcal{A}} \|A\sigma'\| \right) \cdot \sup_{A \in \mathcal{A}} \sum_{r=1}^{\infty} 2^{r/2} \|\Delta_r A\| \\ &\lesssim \left( \mathbb{E}_{\sigma'} \sup_{A \in \mathcal{A}} \|A\sigma'\| \right) \cdot \sup_{A \in \mathcal{A}} \sum_{r=1}^{\infty} 2^{r/2} \cdot d_{\ell_2 \rightarrow \ell_2}(A, T_r), \end{aligned}$$

since  $\|\Delta_r A\| \leq d_{\ell_2 \rightarrow \ell_2}(A, T_{r-1}) + d_{\ell_2 \rightarrow \ell_2}(A, T_r)$  via the triangle inequality. Choosing admissible  $T_0 \subseteq T_1 \subseteq \dots \subseteq T$  to minimize the above expression,

$$E \lesssim \rho_F(\mathcal{A}) \cdot \rho_{\ell_2 \rightarrow 2}(\mathcal{A}) + \gamma_2(\mathcal{A}, \|\cdot\|) \cdot \mathbb{E}_{\sigma'} \sup_{A \in \mathcal{A}} \|A\sigma'\|.$$

Now observe

$$\begin{aligned} \mathbb{E}_{\sigma'} \left( \sup_{A \in \mathcal{A}} \|A\sigma'\| \right) &\leq \left( \mathbb{E}_{\sigma'} \sup_{A \in \mathcal{A}} \|A\sigma'\|^2 \right)^{1/2} \\ &\leq \left( \mathbb{E}_{\sigma'} \left( \sup_{A \in \mathcal{A}} \left| \|A\sigma'\|^2 - \mathbb{E}_{\sigma'} \|A\sigma'\|^2 \right| + \mathbb{E}_{\sigma'} \|A\sigma'\|^2 \right) \right)^{1/2} \\ &= \left( \mathbb{E}_{\sigma'} \sup_{A \in \mathcal{A}} \left( \left| \|A\sigma'\|^2 - \mathbb{E}_{\sigma'} \|A\sigma'\|^2 \right| + \|A\|_F^2 \right) \right)^{1/2} \end{aligned}$$

$$\leq \sqrt{E} + \rho_F(\mathcal{A})$$

Thus in summary,

$$E \lesssim \rho_F(\mathcal{A}) \cdot \rho_{\ell_2 \rightarrow \ell_2}(\mathcal{A}) + \gamma_2(\mathcal{A}, \|\cdot\|) \cdot (\sqrt{E} + \rho_F(\mathcal{A})).$$

This implies  $E$  is at most the square of the larger root of the associated quadratic equation, which gives the theorem.  $\square$

Now we use the KMR theorem (together with majorizing measures) to recover [Theorem 8.2.1](#) (see also [\[KM05, MPTJ07, Dir16\]](#)). We again only discuss the Rademacher case.

**Theorem 8.2.3.** *Let  $T \subset \mathbb{R}^d$  be a set of vectors each of unit norm, and let  $\varepsilon \in (0, 1/2)$  be arbitrary. Let  $\Pi \in \mathbb{R}^{m \times d}$  be such that  $\Pi_{i,j} = \sigma_{i,j}/\sqrt{m}$  for independent Rademacher  $\sigma_{i,j}$ , and where  $m = \Omega((\gamma_2^2(T, \|\cdot\|) + 1)/\varepsilon^2)$ . Then*

$$\mathbb{E} \sup_{x \in T} \left| \|\Pi x\|^2 - 1 \right| < \varepsilon.$$

*Proof.* Similarly to [Eq. \(5.4\)](#), for  $x \in T$  let  $A_x$  denote the  $m \times md$  matrix defined as follows:

$$A_x = \frac{1}{\sqrt{m}} \cdot \begin{bmatrix} x_1 & \cdots & x_d & 0 & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & 0 \\ 0 & \cdots & 0 & x_1 & \cdots & x_d & 0 & \cdots & \cdots & \cdots & \cdots & 0 \\ \vdots & \vdots & \vdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & 0 & x_1 & \cdots & x_d \end{bmatrix}.$$

Then  $\|\Pi x\|^2 = \|A_x \sigma\|^2$ , so letting  $\mathcal{A} = \{A_x : x \in T\}$ ,

$$\mathbb{E} \sup_{x \in T} \left| \|\Pi x\|^2 - 1 \right| = \mathbb{E} \sup_{A \in \mathcal{A}} \left| \|A\sigma\|^2 - \mathbb{E} \|A\sigma\|^2 \right|.$$

We have  $\rho_F(\mathcal{A}) = 1$ . Also  $A_x^\top A_x$  is a block-diagonal matrix, with  $m$  blocks each equal to  $xx^\top/m$ , and thus the singular values of  $A_x$  are 0 and  $\|x\|/\sqrt{m}$ , implying  $\rho_{\ell_2 \rightarrow \ell_2}(\mathcal{A}) = 1/\sqrt{m}$ . Similarly, since  $A_x - A_y = A_{x-y}$ , for any vectors  $x, y$  we have  $\|A_x - A_y\| = \|x - y\|$ , and thus  $\gamma_2(\mathcal{A}, \|\cdot\|) = \gamma_2(T, \|\cdot\|)/\sqrt{m}$ . Thus by the KMR theorem we have

$$\mathbb{E} \sup_{x \in T} \left| \|\Pi x\|^2 - 1 \right| \lesssim \frac{\gamma_2^2(T, \|\cdot\|)}{m} + \frac{\gamma_2(T, \|\cdot\|)}{\sqrt{m}} + \frac{1}{\sqrt{m}},$$

which is at most  $\varepsilon$  for  $m$  as in the theorem statement.  $\square$

Gordon's theorem was actually stated differently in [\[Gor88\]](#) in two ways: (1) Gordon actually only analyzed the case of  $\Pi$  having i.i.d. gaussian entries, and (2) the  $\gamma_2(T, \|\cdot\|)$  terms in the theorem statement were written as the mean width  $w(T)$ . For (1), the extension to arbitrary subgaussian random variables was shown first in [\[KM05\]](#). Note the KMR theorem only bounds an expectation; thus if one wants to argue that the random variable in question is large with probability at most  $\delta$ , the most obvious way is Markov, which would introduce JL a poor  $1/\delta^2$  dependence in  $m$ . One could remedy this by doing Markov on the  $p$ th moment; the tightest known  $p$ -norm bound is given in [\[Dir13, Theorem 6.5\]](#) (see also [\[Dir16, Theorem 4.8\]](#)).

For (2), Gordon actually wrote his paper before  $\gamma_2$  was even defined! The definition of  $\gamma_2$  given here is due to Talagrand, who also showed that for all sets of vectors  $T \subset \mathbb{R}^n$ ,  $w(T) \simeq \gamma_2(T, \|\cdot\|)$  [\[Tal14\]](#) (this is known as the ‘‘majorizing measures’’ theorem). In fact the upper bound  $w(T) \lesssim \gamma_2(T, \|\cdot\|)$  was shown by Fernique [\[Fer75\]](#) (although  $\gamma_2$  was not defined at that point; Talagrand later recast this upper bound in terms of his newly defined  $\gamma_2$ -functional).

[Theorem 8.2.1](#) is thus a corollary of [Theorem 8.2.3](#) and the majorizing measures theorem.  $\square$

### 8.3 Heavy hitters: the BPTree

In this section, we show a result of [BCI<sup>+</sup>17] (building upon [BCIW16]) that leverages tools from the study of suprema of stochastic processes (specifically Dudley’s inequality) to design an  $\ell_2$  heavy hitters algorithm in insertion-only streams using only  $O(k \log k)$  words of memory. Contrast this with the CountSketch of Subsection 4.1.2, which uses  $O(k \log n)$  words of memory for the same task (where we have an insertion-only stream over the universe  $[n]$  updating a frequency vector  $z \in \mathbb{R}^n$ ). Similarly to Subsection 4.1.2, in this section we focus on the tail heavy hitters guarantee. That is, our data structure must respond to a query with a set  $L \subseteq [n]$  that with probability at least  $2/3$  satisfies both (1)  $|L| = O(k)$ , and (2)  $z_i > \frac{1}{\sqrt{k}} \|z_{\text{tail}(k)}\|_2 \implies i \in L$ . In this section we refer to an index  $i$  satisfying (2) as a heavy hitter. Note the number of heavy hitters is at most  $2k - 1$  (namely the  $k$  elements in the support of  $z_{\text{head}(k)}$ , as well as the at most  $k - 1$  elements from the tail satisfying (2)). The correctness conditions then essentially say  $L$  should not be more than a constant factor times bigger than its maximum possible size, and there should also be no false negatives.

Before we describe the BPTree and its analysis, we first prove a modified form of Dudley’s inequality that holds even under limited independence. We use the same definition of  $\Delta_k x$  as in Section 8.1. Note the right hand side of the lemma conclusion is bounded by  $(\mathcal{N}(T, \ell_2, 1/2^k) \cdot \mathcal{N}(T, \ell_2, 1/2^{k-1}))^{1/p} \leq (\mathcal{N}(T, \ell_2, 1/2^k))^{2/p}$ , but we write the more precise form below as it can help us reduce the choice of  $p$  in some applications by a factor of 2 (as we shall see later in this section).

**Lemma 8.3.1.** *Let  $\sigma_1, \dots, \sigma_n$  be  $p$ -wise independent Rademachers for  $p$  an even integer. Then*

$$\mathbb{E} \sup_{x \in T} |\langle \sigma, x \rangle| \lesssim \sqrt{p} \cdot \sum_{k=1}^{\infty} \frac{1}{2^k} \cdot (|\{x' : \exists x \in T \ x' = \Delta_k x\}|)^{1/p}$$

*Proof.* Recall that by applying Lemma 1.1.15 to Khintchine’s inequality (Theorem 1.1.7), for independent  $\sigma_i$  we have that

$$\forall q \geq 1, \|\langle \sigma, x \rangle\|_q \leq \sqrt{q} \|x\|_2$$

Note also that since  $p$  is even,  $\|\langle \sigma, x \rangle\|_p^p = \mathbb{E} \langle \sigma, x \rangle^p$  (that is, the absolute value may be dropped), and the right hand side is equal to

$$\sum_{i_1, \dots, i_p} \left( \prod_{j=1}^p x_{i_j} \right) \left( \mathbb{E} \prod_{j=1}^p \sigma_{i_j} \right),$$

which is fully determined by  $p$ -wise independence. Our proof is then quite similar to that of Eq. (8.1), except that we are only allowed to use tail bounds implied by the moment version of Khintchine’s inequality up to the  $p$ th moment. Observe then that even under  $p$ -wise independence,

$$\mathbb{P}_{\sigma}(|\langle \sigma, x \rangle| > \lambda) < \left( \frac{\sqrt{p} \|x\|_2}{\lambda} \right)^p \tag{8.3}$$

Now first we show what this implies when using the union bound. Suppose  $T \subseteq B_{\ell_2}^n$ . Then

$$\begin{aligned} \mathbb{E} \sup_{x \in T} |\langle \sigma, x \rangle| &= \int_0^{\infty} \mathbb{P}(\sup_{x \in T} |\langle \sigma, x \rangle| > t) dt \\ &= \underbrace{\int_0^{t^*} \mathbb{P}(\sup_{x \in T} |\langle \sigma, x \rangle| > t) dt}_{\leq 1} + \int_{t^*}^{\infty} \mathbb{P}(\sup_{x \in T} |\langle \sigma, x \rangle| > t) dt \end{aligned}$$



$$\begin{aligned}
&\leq t^* + \sum_{x \in T} \int_{t^*}^{\infty} \mathbb{P}(|\langle \sigma, x \rangle| > t) dt \\
&\leq t^* + \frac{|T|}{p-1} \cdot \left( \frac{\sqrt{p}}{t^*} \right)^{p-1} && \text{(by Eq. (8.3))} \\
&\lesssim \sqrt{p} \cdot |T|^{1/p} && (8.4)
\end{aligned}$$

where in the final equality we chose  $t^*$  to make the two summands equal.

Then using the same definitions for  $\pi_k x, \Delta_k x, S_k$  as in [Section 8.1](#),

$$\begin{aligned}
\mathbb{E} \sup_{x \in T} |\langle \sigma, x \rangle| &\leq \sum_{k=1}^{\infty} \mathbb{E} \sup_{x \in T} |\langle \sigma, \Delta_k x \rangle| \\
&\leq \sum_{k=1}^{\infty} \left( \sup_{x \in T} \|\Delta_k x\|_2 \right) \cdot \mathbb{E} \sup_{x \in T} \left| \langle \sigma, \frac{\Delta_k x}{\|\Delta_k x\|_2} \rangle \right| \\
&\lesssim \sqrt{p} \cdot \sum_{k=1}^{\infty} \left( \sup_{x \in T} (\|\pi_k x\|_2 + \|\pi_{k-1} x\|_2) \right) \cdot (|\{x' : \exists x \in T \ x' = \Delta_k x\}|)^{1/p} \\
&\lesssim \sqrt{p} \cdot \sum_{k=1}^{\infty} \frac{1}{2^k} \cdot (|\{x' : \exists x \in T \ x' = \Delta_k x\}|)^{1/p} && (8.5)
\end{aligned}$$

□

We now state a lemma that will be useful for the BPTree.

**Lemma 8.3.2.** *Let  $y$  be a frequency vector updated in an insertion-only stream, and  $y^{(t)}$  denote this vector at time  $t$  (i.e. after the first  $t$  updates). Let  $\sigma_1, \dots, \sigma_n$  be 4-wise independent  $\pm 1$ . Then*

$$\mathbb{E} \max_{1 \leq t \leq T} |\langle \sigma, y^{(t)} \rangle| \lesssim \|y^{(T)}\|_2$$

*Proof.* Define  $v^{(t)} := y^{(t)} / \|y^{(T)}\|_2 \in B_{\ell_2^2}$  and  $V := \{v^{(t)}\}_{t=1}^T$ . We apply [Lemma 8.3.1](#) to  $V$ . For this, we first understand the structure of  $\varepsilon$ -nets of  $V$ . We give a greedy construction: our  $\varepsilon$ -net is  $\{v^{(t_0)}, \dots, v^{(t_r)}\}$  where  $t_0 = 0$  and for  $j \geq 1$  we define  $t_j$  to be the smallest time  $t > t_{j-1}$  such that  $\|v^{(t)} - v^{(t_{j-1})}\|_2 > \varepsilon$  (if such  $t$  exists). We first observe that  $r < 1/\varepsilon^2$ . To see this, define  $w_j := v^{(t_j)} - v^{(t_{j-1})}$ . Then

$$\begin{aligned}
1 &\geq \|v^{(t_r)}\|_2^2 \\
&= \left\| \sum_{j=1}^r w_j \right\|_2^2 \\
&= \sum_{j=1}^r \|w_j\|_2^2 + \sum_{j \neq j'} \langle w_j, w_{j'} \rangle \\
&\geq \sum_{j=1}^r \|w_j\|_2^2 && (w_j, w_{j'} \text{ have nonnegative coordinates}) \\
&> r\varepsilon^2
\end{aligned}$$

so that  $r < 1/\varepsilon^2$ . We thus let the  $S_k$  be a sequence of  $2^{-k}$ -nets where  $|S_k| < 1/\varepsilon^2 + 1$ . Next observe that for any  $k \geq 1$ , using the above-described nets, the number of possibilities for  $\Delta_k v$  is at

most  $2|S_k|$  (as opposed to the trivial bound  $|S_k| \cdot |S_{k-1}|$ ). This is because the net elements can be arranged linearly (sorted by time), and given that  $v^{(t')} = \pi_k v^{(t)}$ , there are only two possibilities for  $\pi_{k-1} v^{(t)}$ : either the closest net point backward in time from  $t'$  in  $S_{k-1}$ , or forward in time. Since  $|S_k| = O(2^{2k})$ , Lemma 8.3.1 yields  $\mathbb{E} \sup_{v \in V} |\langle \sigma, v \rangle| \lesssim \sum_{k=1}^{\infty} \frac{1}{2^k} \cdot (2^{2k})^{1/4} = O(1)$ , as desired.  $\square$

Narayanan showed that 4-wise independence is not only sufficient, but also necessary for the conclusion of Lemma 8.3.2 to hold [Nar19]. That is, there are 3-wise independent distributions over  $\sigma_1, \dots, \sigma_n$  where the conclusion of the lemma does not hold.

### The algorithm and analysis.

We now describe the BPTree. We first show a reduction from the heavy hitters problem to the *super heavy hitter* problem. In particular, we show that a solution to the latter with failure probability  $1/10$  in space  $S$  words implies a solution to the former with failure probability  $1/k^c$  for any constant  $c$  with space  $O(Sk \log k)$  [BCIW16]. We then show that the super heavy hitter problem can be solved with  $S = O(1)$  memory.

**Definition 8.3.3.** We say an index  $i \in [n]$  is *super-heavy* if  $\sum_i z_i^2 > 1000 \cdot \sum_{j \neq i} z_j^2$ .

The definition implies that there can be at most one super-heavy element in a stream. As an algorithmic problem, we would like a solution which (1) is allowed to behave arbitrarily if the stream does not have any super heavy, and (2) outputs the unique super-heavy element with probability at least  $9/10$ , if it exists.

**Reduction to super heavy hitter.** We reduce the heavy hitters problem to super heavy via hashing. Suppose we have a solution to the super heavy hitter problem with failure probability  $1/10$ . Set  $q = Ck$ ,  $r = C \log k$  for a large constant  $C > 0$  and initialize  $qr$  copies of this data structure  $B_{t,j}$  for  $t \in [r], j \in [q]$ . Also, draw hash functions  $h_1, \dots, h_r : [n] \rightarrow [q]$  independently from a pairwise independent hash family.

When we see  $i$  in the stream, we feed  $i$  to  $B_{t,h_t(j)}$  for each  $t \in [r]$ . When queried for the list  $L$  of heavy hitters, we return  $L = \{i \in [n] : B_{t,j}.\text{query}() \text{ returns } i \text{ for at least } r/2 \text{ values of } (t, j) \in [r] \times [q]\}$ . Thus if each  $B_{t,j}$  uses space  $S$ , has update time  $t_u$  and query time  $t_q$ , our final space is  $O(Sk \log k)$ , our update time is  $O(t_u r) = O(t_u \log k)$ , and our query time is  $O(k \log^2 k t_q)$  (or  $O(k \log k t_q)$  expected time if using a dictionary).

**Theorem 8.3.4.** *The failure probability of the above reduction is at most  $1/k^c$ , where  $c$  can be made an arbitrarily large constant by increasing  $C$ .*

*Proof.* By construction, the output list  $L$  is guaranteed to have size at most  $2q = O(k)$ . Thus we must only bound the probability of a false negative.

Let  $H \subset [n]$  be the set of all heavy hitters. Consider a specific heavy hitter  $i \in H$ . We will show that the probability that  $i \notin L$  is at most  $1/(2k^{c+1})$ , so that the probability of any heavy hitter failing to be in  $L$  is at most  $|H|/(2k^{c+1}) < 1/k^c$  by a union bound. Fix a particular  $t \in [r]$ . Define the event  $\mathcal{E}_1(t)$  that no  $j \in \text{head}(k)$  collides with  $i$  under  $h_t$ , i.e. the event that  $h_t^{-1}(h_t(i)) \cap \text{head}(k) = \{i\}$ . Also define the event  $\mathcal{E}_2(t)$  that  $\sum_{j \in \text{tail}(k): h_t(j) = h_t(i)} z_j^2 \geq \|z_{\text{tail}(k)}\|_2^2 / (1000k)$ . Note if both  $\mathcal{E}_1(t), \mathcal{E}_2(t)$  occur, then  $i$  is super heavy in the stream processed by  $B_{t,h_t(i)}$ .  $\mathbb{P}(\neg \mathcal{E}_1(t)) \leq k/q \leq 1/10$  for  $C$  sufficiently large. To see this, one can define indicator random variables  $X_j(t)$  that are 1 iff  $h_t(j) = h_t(i)$ . Then

$$\mathbb{E} \sum_{j \in \text{head}(k) \setminus \{i\}} X_j(t) = \sum_{j \in \text{head}(k) \setminus \{i\}} \mathbb{E} X_j(t) = \frac{|\text{head}(k)|}{q} = \frac{k}{q},$$

since  $\mathbb{E} X_j(t) = 1/q$  by pairwise independence of  $h_t$ . Then  $\neg \mathcal{E}_1(t)$  holds iff the above summation is at least 1, which happens with probability less than  $k/q$  by Markov's inequality. Similarly for  $\mathcal{E}_2$ ,

$$\mathbb{E} \sum_{j \in \text{tail}(k)} X_j(t) \cdot z_j^2 = \frac{\|z_{\text{tail}(k)}\|_2^2}{q},$$

and thus the probability that the above sum is at least  $\|z_{\text{tail}(k)}\|_2^2/(1000k)$  is at most  $1000k/q$ , which is also at most  $1/10$  for  $C$  sufficiently large. Thus  $i$  is super heavy to  $B_{t, h_t(i)}$  with probability at least  $4/5$ , implying  $i$  is output by  $B_{t, h_t(i)}$  with probability at least  $(4/5) \cdot (9/10) > 7/10$ .

Now observe that the expected number of  $t \in [r]$  for which  $B_{t, h_t(i)}$  outputs  $i$  is  $7r/10$ . Thus the probability that fewer than  $r/2$  such data structures output  $i$  is at most  $\exp(-\Omega(r))$  by the Chernoff bound. This is at most  $1/(2k^{c+1})$  for  $C$  a sufficiently large constant.  $\square$

**An algorithm for super heavy hitter.** We now describe an algorithm for super heavy hitter using  $O(1)$  words of memory.

Let the stream length be  $\ell$  and let  $T = \|z^{(\ell)}\|_2$  be the  $\ell_2$  norm of the frequency vector at the end of the stream. We assume we know  $T$  exactly (an assumption we will remove later), and we also assume the stream has some super heavy item  $i^* \in [n]$  since otherwise the algorithm is allowed to behave arbitrarily. We further assume that  $z_{i^*}^2 \geq C \sum_{i \neq i^*} z_i^2$  for some large constant  $C$  that is implicit in the proof (it will be larger than 1000). Though the reduction to the super heavy problem took  $C = 1000$ , the proof above goes through essentially verbatim for any large constant. We select  $h : [N] \rightarrow [N]$  from a 2-wise independent family, then we attempt to learn the  $\log_2 N$  bits of  $h(i^*)$  one at a time, where  $N$  is  $Cn^2$ . Then once we learn  $h(i^*)$ , we simply return the next index  $i$  in the stream such that  $h(i) = h(i^*)$  (we will guarantee that we learn  $i^*$  before its last occurrence). The stream may end before we learn all  $\log N$  bits (e.g. it may be that  $N$  is much bigger than  $T$ ); this is discussed in the next paragraph.

Write  $h(i^*)$  in binary as  $\sum_{j=0}^t b_j 2^j$ , where  $t = \lfloor \log_2 N \rfloor$ . We aim to learn the  $b_j$  one bit at a time, iteratively. Suppose we have learned  $b_t, b_{t-1}, \dots, b_{t-r+1}$  so far for  $r \geq 0$  and would now like to learn  $b_{t-r}$  (initially  $r = 0$ ). We thus know that  $h(i^*) \in \mathcal{I}$ , where  $\mathcal{I}$  is the interval  $[B \cdot 2^{t-r+1}, (B+1) \cdot 2^{t-r+1}) \subseteq \{0, \dots, N-1\}$  and  $B$  is the  $r$ -bit number  $b_t b_{t-1} \dots b_{t-r+1}$ . We draw  $\sigma \in \{-1, 1\}^n$  from a 4-wise independent family and initialize two counters  $C_0, C_1$  each to zero. When we see  $i$  in the stream, if  $h(i) \notin \mathcal{I}$  then we ignore  $i$  and continue ( $i$  cannot possibly be  $i^*$ , so it is not worth processing). Otherwise, we look at the  $(t-r)$ th bit of  $h(i)$ : if it is  $b$ , then we increment  $C_b$  by  $\sigma_i$ . We continue processing the stream in this way until the first time that  $|C_b| \geq \lceil (T/11) \cdot (2/3)^{r/2} \rceil$  for some  $b \in \{0, 1\}$ . We then declare that the  $(t-r)$ th bit of  $h(i^*)$  is  $b$ , increment  $r$ , then iterate. We halt after learning all  $\lfloor \log_2 N \rfloor + 1$  bits of  $h(i^*)$ , i.e. once we've learned bit  $b_0$ , then we return the next stream element  $i$  such that  $h(i)$  matches the bits we have learned so far (thus we may not have learned all  $\log N$  bits of  $h(i^*)$ , but we have learned enough that a unique index in the stream matches these bits). This concludes our description of the algorithm.

We now proceed to analysis. Let us write the  $(t-r)$ th bit of  $h(i^*)$  as  $b_r^* \in \{0, 1\}$ . Let  $X'_i(r)$  be the indicator random variable for the event that the  $j$ th bit of  $h(i)$  agrees with the  $j$ th bit of  $h(i^*)$  for  $j = t, \dots, t-r$  and  $X''_i(r)$  be the indicator for the event that the  $j$ th bits agree for  $j = t, \dots, t-r+1$  but then the  $(t-r)$ th bit differs. We also define  $X'_{i^*}(r) = X''_{i^*}(r) = 0$ . Then we can define vectors  $z'(r)$  and  $z''(r)$ , where  $z'(r)_i = X'_i(r) z_i$ ,  $z''(r)_i = X''_i(r) z_i$ , so that

$$C_{1-b_r^*} = \underbrace{\langle \sigma, z''(r) \rangle}_{\alpha''(r)}, \quad C_{b_r^*} = \sigma_{i^*} z_{i^*} + \underbrace{\langle \sigma, z'(r) \rangle}_{\alpha'(r)}$$

We also let  $z'^{(t)}$  denote  $z'$  after processing the  $t$ th update (“time  $t$ ”), and  $z''$  similarly, and we use  $C_{b_r^*}^{(t)}$  to denote  $\langle \sigma, z'^{(t)} \rangle$  (and similarly for  $C_{1-b_r^*}^{(t)}$ ). We also similarly define  $\alpha'(r)^{(t)}, \alpha''(r)^{(t)}$ .

**Lemma 8.3.5.** *For  $C$  sufficiently large, with probability at least  $9/10$*

$$\forall r, t \geq 0, |\alpha'(r)^{(t)}|, |\alpha''(r)^{(t)}| \leq \frac{z_{i^*}}{100} \cdot \left(\frac{2}{3}\right)^{r/2}. \quad (8.6)$$

*Proof.* By Lemma 8.3.2, there is a constant  $c > 0$  such that for any fixed  $r \geq 0$  we have

$$\begin{aligned} \mathbb{E} \max_{\sigma, h} \max_{1 \leq t \leq \ell} |\alpha'(r)^{(t)}| &\leq c \cdot \mathbb{E}_h \|z'(r)^{(\ell)}\|_2 && \text{(Lemma 8.3.2)} \\ &\leq c \cdot (\mathbb{E}_h \|z'(r)^{(\ell)}\|_2^2)^{1/2} && \text{(Cauchy-Schwarz)} \\ &= c \cdot \left( \sum_{i \neq i^*} (\mathbb{E} X_i'(r)) z_i^2 \right)^{1/2} \\ &= \frac{c}{2^{(r+1)/2}} \cdot \left( \sum_{i \neq i^*} z_i^2 \right)^{1/2} \\ &< \left( \frac{c}{\sqrt{C}} \right) \cdot \frac{z_{i^*}}{2^{(r+1)/2}} && \text{(definition of super heaviness)} \end{aligned}$$

Thus by Markov’s inequality,

$$\mathbb{P}_h \left( \max_{1 \leq t \leq \ell} |\alpha'(r)^{(t)}| > \frac{75\sqrt{2}cz_{i^*}}{\sqrt{C}} \cdot \left(\frac{2}{3}\right)^{r/2} \right) < \frac{1}{150} \cdot \left(\frac{3}{4}\right)^{r/2}.$$

Therefore by a union bound over all  $r \geq 0$ ,

$$\mathbb{P}_h \left( \exists r, t \geq 0, |\alpha'(r)^{(t)}| > \frac{75\sqrt{2}cz_{i^*}}{\sqrt{C}} \cdot \left(\frac{2}{3}\right)^{r/2} \right) < \sum_{r=0}^{\infty} \frac{1}{150} \cdot \left(\frac{3}{4}\right)^{r/2} < \frac{1}{20}.$$

An identical argument applies to  $|\alpha''(r)^{(t)}|$ . The lemma statement the holds by a union bound over both  $|\alpha'(r)^{(t)}|$  and  $|\alpha''(r)^{(t)}|$  for all  $r, t \geq 0$  and setting  $C$  sufficiently large ( $C \geq 2 \cdot 75^2 100^2 c^2$ ).  $\square$

We condition on the event of Eq. (8.6). Then since  $|\alpha''(r)| = |C_{1-b_r^*}| \leq (z_i^*/100)(2/3)^{r/2}$  for each  $r$ , we will never incorrectly learn  $b_r^*$  for any  $r$ . Furthermore since  $|\alpha'(r)|$  is similarly small, we will learn  $b_r^*$  in the next at most  $\lceil (T/11) \cdot (2/3)^{r/2} \rceil + (z_i^*/100)(2/3)^{r/2} \leq \lceil z_i^*/10(2/3)^{r/2} \rceil$  occurrences of  $i^*$  in the stream after iteration  $r$  begins.

Note:  $i^*$  is the unique index still in  $\mathcal{I}$  after the first  $O(\log \min\{n, T\}) + z_i^*/9$  occurrences of  $i^*$  in the stream, at which point we can simply return the next index  $i$  in  $\mathcal{I}$  (i.e. the next index which matches the bits we’ve learned so far).

**Removing the assumption of knowing  $\|z^{(\ell)}\|_2$  exactly.** We do not provide the full details here, but only just a sketch. The first observation is that we do not need to know  $\|z^{(\ell)}\|_2$  exactly, but rather it suffices to know a value  $T$  within a factor of two of this quantity. Since we do not know this either, we guess  $T = 1, 2, \dots, 2^{10}$  in parallel. We simultaneously run the AMS sketch in parallel using 8-wise independence; it can be shown that the AMS sketch with this increased independence

and  $O(1/\alpha^2)$  rows has the property that even if you query after every single update, with large probability there will never be any time  $t$  such that the AMS output is not  $\|z^{(t)}\|_2 \pm \alpha\|z^{(\ell)}\|_2$  (see [BCI<sup>+</sup>17]); we apply this with  $\alpha$  a small constant. We call this guarantee  $\ell_2$ -tracking, since we can track  $\|z^{(t)}\|_2$  at all times  $t$ . Note that naively one would instantiate the failure probability of the AMS sketch to be  $\ll 1/\ell$  then union bound over all time steps, but that would require memory  $O(\alpha^{-2} \log \ell)$  words, which we avoid. Now, let us say at some time step we are running the algorithm in parallel above for the guesses of  $T$  being  $2^j, \dots, 2^{j+10}$ . When the tracker says that the guess  $2^j$  is too small, we simply kill that parallel instantiation, reclaim the space, then boot up a new parallel instantiation which guesses  $T = 2^{j+11}$ . This new instantiation has missed the prefix of the stream up until this point, but in the case that  $\|z^{(\ell)}\|_2 \approx 2^{j+11}$ , the prefix missed was insignificant (since it only contained less than 0.1% of the occurrences of  $i^*$ ).



# Bibliography

- [AC09] Nir Ailon and Bernard Chazelle. The fast Johnson–Lindenstrauss transform and approximate nearest neighbors. *SIAM J. Comput.*, 39(1):302–322, 2009. [66](#), [68](#), [69](#), [70](#), [78](#)
- [Ach03] Dimitris Achlioptas. Database-friendly random projections: Johnson-lindenstrauss with binary coins. *J. Comput. Syst. Sci.*, 66(4):671–687, 2003. [66](#)
- [ACH<sup>+</sup>13] Pankaj K. Agarwal, Graham Cormode, Zengfeng Huang, Jeff M. Phillips, Zhewei Wei, and Ke Yi. Mergeable summaries. *ACM Trans. Database Syst.*, 38(4):26:1–26:28, 2013. [28](#), [29](#)
- [AGM12] Kook Jin Ahn, Sudipto Guha, and Andrew McGregor. Analyzing graph structure via linear measurements. In *Proceedings of the Twenty-Third Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 459–467, 2012. [49](#)
- [AHLW16] Yuqing Ai, Wei Hu, Yi Li, and David P. Woodruff. New characterizations in turnstile streams with applications. In *Proceedings of the 31st Conference on Computational Complexity (CCC)*, pages 20:1–20:22, 2016. [43](#)
- [AK17] Noga Alon and Bo’az Klartag. Optimal compression of approximate inner products and dimension reduction. In *Proceedings of the 58th IEEE Annual Symposium on Foundations of Computer Science (FOCS)*, pages 639–650, 2017. [62](#)
- [AL09] Nir Ailon and Edo Liberty. Fast dimension reduction using rademacher series on dual BCH codes. *Discret. Comput. Geom.*, 42(4):615–630, 2009. [66](#)
- [AL13] Nir Ailon and Edo Liberty. An almost optimal unrestricted fast Johnson-Lindenstrauss transform. *ACM Trans. Algorithms*, 9(3):21:1–21:12, 2013. [66](#), [70](#)
- [Alo03] Noga Alon. Problems and results in extremal combinatorics–I. *Discrete Mathematics*, 273(1-3):31–53, 2003. [61](#), [65](#)
- [AMS99] Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. *J. Comput. Syst. Sci.*, 58(1):137–147, 1999. [20](#), [31](#), [53](#)
- [AW20] Josh Alman and Virginia Vassilevska Williams. A refined laser method and faster matrix multiplication. *CoRR*, abs/2010.05846, 2020. [75](#)
- [BCI<sup>+</sup>17] Vladimir Braverman, Stephen R. Chestnut, Nikita Ivkin, Jelani Nelson, Zhengyu Wang, and David P. Woodruff. BPTree: An  $\ell_2$  heavy hitters algorithm using constant memory. In *Proceedings of the 36th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (PODS)*, pages 361–376, 2017. [112](#), [117](#)

- [BCIW16] Vladimir Braverman, Stephen R. Chestnut, Nikita Ivkin, and David P. Woodruff. Beating counts sketch for heavy hitters in insertion streams. In *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 740–753, 2016. [112](#), [114](#)
- [BD09] Thomas Blumensath and Mike E. Davies. A simple, efficient and near optimal algorithm for compressed sensing. In *ICASSP*, 2009. [101](#)
- [BDF<sup>+</sup>11] Jean Bourgain, Stephen Dilworth, Kevin Ford, Sergei Konyagin, and Denka Kutzarova. Explicit constructions of rip matrices and related problems. *Duke Mathematical Journal*, pages 145–185, 2011. [100](#)
- [BDN15] Jean Bourgain, Sjoerd Dirksen, and Jelani Nelson. Toward a unified theory of sparse dimensionality reduction in Euclidean space. *Geometric and Functional Analysis (GAFA)*, to appear, 2015. Preliminary version in STOC 2015. [82](#)
- [BJKS04] Ziv Bar-Yossef, T. S. Jayram, Ravi Kumar, and D. Sivakumar. An information statistics approach to data stream and communication complexity. *J. Comput. Syst. Sci.*, 68(4):702–732, 2004. [38](#), [39](#), [53](#)
- [Bła20] Jarosław Błasiok. Optimal streaming and tracking distinct elements with high probability. *ACM Trans. Algorithms*, 16(1):3:1–3:28, 2020. [20](#)
- [BOR10] Vladimir Braverman, Rafail Ostrovsky, and Yuval Rabani. Rademacher chaos, random Eulerian graphs and the sparse Johnson-Lindenstrauss transform. *CoRR*, abs/1011.2590, 2010. [66](#)
- [Bou14] Jean Bourgain. An improved estimate in the restricted isometry problem. *Geometric Aspects of Functional Analysis*, Lecture Notes in Mathematics Volume 2116:65–70, 2014. [71](#), [100](#)
- [BR94] Mihir Bellare and John Rompel. Randomness-efficient oblivious sampling. In *Proceedings of the 35th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 276–287, 1994. [51](#)
- [Bro93] Andrej Brodnik. Computation of the least significant set bit. In *ERK*, 1993. [23](#)
- [BSS12] Joshua D. Batson, Daniel A. Spielman, and Nikhil Srivastava. Twice-ramanujan sparsifiers. *SIAM J. Comput.*, 41(6):1704–1721, 2012. [82](#)
- [BYJK<sup>+</sup>02] Ziv Bar-Yossef, T.S. Jayram, R. Kumar, D. Sivakumar, and Luca Trevisan. Counting distinct elements in a data stream. In *RANDOM*, pages 1–10, 2002. [19](#)
- [BZMD15] Christos Boutsidis, Anastasios Zouzias, Michael W. Mahoney, and Petros Drineas. Randomized dimensionality reduction for k-means clustering. *IEEE Trans. Inf. Theory*, 61(2):1045–1062, 2015. [90](#)
- [Can08] Emmanuel Candès. The restricted isometry property and its implications for compressed sensing. *Comptes Rendus Mathématique*, 346(9-10):589–592, 2008. [97](#)
- [CCF04] Moses Charikar, Kevin C. Chen, and Martin Farach-Colton. Finding frequent items in data streams. *Theor. Comput. Sci.*, 312(1):3–15, 2004. [47](#)



- [CEM<sup>+</sup>15] Michael B. Cohen, Sam Elder, Cameron Musco, Christopher Musco, and Madalina Persu. Dimensionality reduction for  $k$ -means clustering and low rank approximation. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing (STOC)*, pages 163–172, 2015. 89, 90, 93
- [CJN18] Michael B. Cohen, T.S. Jayram, and Jelani Nelson. Simple analyses of the Sparse Johnson-Lindenstrauss transform. In *Proceedings of the 1st Annual Symposium on Simplicity in Algorithms (SOSA)*, 2018. 66, 67
- [CK16] Amit Chakrabarti and Sagar Kale. Strong fooling sets for multi-player communication with applications to deterministic estimation of stream statistics. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 41–50, 2016. 34
- [CKL<sup>+</sup>20] Graham Cormode, Zohar S. Karnin, Edo Liberty, Justin Thaler, and Pavel Veselý. Relative error streaming quantiles. *CoRR*, abs/2004.01668, 2020. 24
- [CM05] Graham Cormode and S. Muthukrishnan. An improved data stream summary: the count-min sketch and its applications. *J. Algorithms*, 55(1):58–75, 2005. 44, 46
- [CMS76] John M. Chambers, Colin L. Mallows, , and B. W. Stuck. A method for simulating stable random variables. *J. Amer. Statist. Assoc.*, 71:340–344, 1976. 56
- [CNW16] Michael B. Cohen, Jelani Nelson, and David P. Woodruff. Optimal approximate matrix product in terms of stable rank. In *Proceedings of the 43rd International Colloquium on Automata, Languages, and Programming (ICALP)*, pages 11:1–11:14, 2016. Full version at <https://arxiv.org/abs/1507.02268v3>. 70, 78, 82
- [Coh16] Michael B. Cohen. Nearly tight oblivious subspace embeddings by trace inequalities. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 278–287, 2016. 82, 83
- [CR12] Amit Chakrabarti and Oded Regev. An optimal lower bound on the communication complexity of gap-hamming-distance. *SIAM J. Comput.*, 41(5):1299–1317, 2012. 40
- [CRT06] Emmanuel J. Candés, Justin K. Romberg, and Terence Tao. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, 2006. 97
- [CT05] Emmaual Candés and Terence Tao. Decoding by linear programming. *IEEE Trans. Inform. Theory*, 51(12):4203–4215, 2005. 70
- [CT06] Emmanuel J. Candès and Terence Tao. Near-optimal signal recovery from random projections: universal encoding strategies? *IEEE Trans. Inform. Theory*, 52:5406–5425, 2006. 71, 100
- [CV20] Graham Cormode and Pavel Veselý. A tight lower bound for comparison-based quantile summaries. In *Proceedings of the 39th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (PODS)*, pages 81–93, 2020. 23, 24
- [CW09] Kenneth L. Clarkson and David P. Woodruff. Numerical linear algebra in the streaming model. In *Proceedings of the 41st Annual ACM Symposium on Theory of Computing (STOC)*, pages 205–214, 2009. 77

- [CW17] Kenneth L. Clarkson and David P. Woodruff. Low rank approximation and regression in input sparsity time. In *J. ACM*, pages 54:1–54:45, 2017. 82, 83, 84
- [CY20] Graham Cormode and Ke Yi. *Small Summaries for Big Data*. Cambridge University Press, 2020. To be published. Draft at <http://dimacs.rutgers.edu/~graham/ssbd.html>. 24, 26
- [Dir13] Sjoerd Dirksen. Tail bounds via generic chaining. *CoRR*, abs/1309.3522v2, 2013. 111
- [Dir16] Sjoerd Dirksen. Dimensionality reduction with subgaussian matrices: a unified theory. *Foundations of Computational Mathematics*, 16(5):1367–1396, 2016. 111
- [DKM06] Petros Drineas, Ravi Kannan, and Michael Mahoney. Fast monte carlo algorithms for matrices I: Approximating matrix multiplication. *SIAM J. Comput*, 36(1):132–157, 2006. 75
- [DKN10] Ilias Diakonikolas, Daniel M. Kane, and Jelani Nelson. Bounded independence fools degree-2 threshold functions. In *51th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 11–20, 2010. 10
- [DKS10] Anirban Dasgupta, Ravi Kumar, and Tamás Sarlós. A sparse Johnson-Lindenstrauss transform. In *Proceedings of the 42nd ACM Symposium on Theory of Computing (STOC)*, pages 341–350, 2010. 66
- [dlPnG99] Victor de la Peña and Evarist Giné. *Decoupling: From dependence to independence*. Probability and its Applications. Springer-Verlag, New York, 1999. 9
- [DMM06] Petros Drineas, Michael W. Mahoney, and S. Muthukrishnan. Sampling algorithms for  $\ell_2$  regression and applications. In *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1127–1136, 2006. 80
- [DMMW12] Petros Drineas, Malik Magdon-Ismail, Michael W. Mahoney, and David P. Woodruff. Fast approximation of matrix coherence and statistical leverage. *J. Mach. Learn. Res.*, 13:3475–3506, 2012. 82
- [Don06] D. Donoho. Compressed sensing. *IEEE Trans. Inform. Theory*, 52(4):1289–1306, 2006. 97
- [FEFGM07] Philippe Flajolet, Éric Fusy, Olivier Gandouet, and Frédéric Meunier. Hyperloglog: the analysis of a near-optimal cardinality estimation algorithm. In *Proceedings of the 2007 International Conference on the Analysis of Algorithms (AoFA)*, 2007. 20
- [Fer75] Xavier Fernique. Régularité des trajectoires des fonctions aléatoires gaussiennes. *Lecture Notes in Math.*, 480:1–96, 1975. 111
- [Fla85] Philippe Flajolet. Approximate counting: A detailed analysis. *BIT Comput. Sci. Sect.*, 25(1):113–134, 1985. 15
- [FM85] Philippe Flajolet and G. Nigel Martin. Probabilistic counting algorithms for data base applications. *J. Comput. Syst. Sci.*, 31(2):182–209, 1985. 16
- [FR13] Simon Foucart and Holger Rauhut. *A Mathematical Introduction to Compressive Sensing*. Birkhäuser, Boston, 2013. 97

- [FW93] Michael L. Fredman and Dan E. Willard. Surpassing the information theoretic bound with fusion trees. *J. Comput. Syst. Sci.*, 47(3):424–436, 1993. [23](#)
- [GJ79] Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. Freeman, San Francisco, CA, 1979. [97](#)
- [GK01] Michael Greenwald and Sanjeev Khanna. Space-efficient online computation of quantile summaries. In *Proceedings of the ACM SIGMOD international conference on Management of data (SIGMOD)*, pages 58–66, 2001. [24](#)
- [GK16] Michael B. Greenwald and Sanjeev Khanna. Quantiles and equi-depth histograms over streams. In Minos N. Garofalakis, Johannes Gehrke, and Rajeev Rastogi, editors, *Data Stream Management - Processing High-Speed Data Streams*, Data-Centric Systems and Applications, pages 45–86. Springer, 2016. [24](#)
- [Gor88] Yehoram Gordon. On Milman’s inequality and random subspaces which escape through a mesh in  $\mathbb{R}^n$ . *Geometric Aspects of Functional Analysis*, pages 84–106, 1988. [108](#), [111](#)
- [GU18] Francois Le Gall and Florent Urrutia. Improved rectangular matrix multiplication using powers of the coppersmith-winograd tensor. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1029–1046, 2018. [75](#)
- [HR16] Ishay Haviv and Oded Regev. The restricted isometry property of subsampled fourier matrices. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 288–297, 2016. [71](#), [100](#)
- [HW71] David Lee Hanson and Farroll Tim Wright. A bound on tail probabilities for quadratic forms in independent random variables. *Ann. Math. Statist.*, 42:1079–1083, 1971. [10](#)
- [ILL<sup>+</sup>19] Nikita Ivkin, Edo Liberty, Kevin J. Lang, Zohar S. Karnin, and Vladimir Braverman. Streaming quantiles algorithms with small space and update time. *CoRR*, abs/1907.00236, 2019. [24](#)
- [Ind06] Piotr Indyk. Stable distributions, pseudorandom generators, embeddings, and data stream computation. *J. ACM*, 53(3):307–323, 2006. [53](#), [54](#)
- [IW05] Piotr Indyk and David P. Woodruff. Optimal approximations of the frequency moments of data streams. In *Proceedings of the 37th Annual ACM Symposium on Theory of Computing (STOC)*, pages 202–208, 2005. [39](#), [53](#)
- [JL84] William B. Johnson and Joram Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. *Contemporary Mathematics*, 26:189–206, 1984. [59](#), [61](#), [66](#)
- [JST11] Hossein Jowhari, Mert Saglam, and Gábor Tardos. Tight bounds for lp samplers, finding duplicates in streams, and related problems. In *Proceedings of the 30th ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS)*, pages 49–58, 2011. [50](#)
- [JW13] T. S. Jayram and David P. Woodruff. Optimal bounds for Johnson-Lindenstrauss transforms and streaming problems with subconstant error. *ACM Transactions on Algorithms*, 9(3):26, 2013. [20](#), [62](#)

- [KLL16] Zohar S. Karnin, Kevin J. Lang, and Edo Liberty. Optimal quantile approximation in streams. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 71–78, 2016. [24](#), [30](#)
- [KM05] Bo’az Klartag and Shahar Mendelson. Empirical processes and random projections. *J. Funct. Anal.*, 225(1):229–245, 2005. [109](#), [111](#)
- [KMN11] Daniel M. Kane, Raghu Meka, and Jelani Nelson. Almost optimal explicit Johnson-Lindenstrauss families. In *RANDOM*, pages 628–639, 2011. [62](#)
- [KMR14] Felix Krahmer, Shahar Mendelson, , and Holger Rauhut. Suprema of chaos processes and the restricted isometry property. *Comm. Pure Appl. Math.*, 67(11):1877–1904, 2014. [109](#)
- [KN97] Eyal Kushilevitz and Noam Nisan. *Communication complexity*. Cambridge University Press, 1997. [37](#)
- [KN10] Daniel M. Kane and Jelani Nelson. A derandomized sparse Johnson-Lindenstrauss transform. *CoRR*, abs/1006.3585, 2010. [66](#)
- [KN14] Daniel M. Kane and Jelani Nelson. Sparser Johnson-Lindenstrauss transforms. *Journal of the ACM*, 61(1):4, 2014. [66](#), [67](#), [78](#)
- [KNP<sup>+</sup>17] Michael Kapralov, Jelani Nelson, Jakub Pachocki, Zhengyu Wang, David P. Woodruff, and Mobin Yahyazadeh. Optimal lower bounds for universal relation, and for samplers and finding duplicates in streams. In *Proceedings of the 58th IEEE Annual Symposium on Foundations of Computer Science (FOCS)*, pages 475–486, 2017. [50](#)
- [KNW10a] Daniel M. Kane, Jelani Nelson, and David P. Woodruff. On the exact space complexity of sketching and streaming small norms. In *Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1161–1178, 2010. [55](#), [58](#)
- [KNW10b] Daniel M. Kane, Jelani Nelson, and David P. Woodruff. An optimal algorithm for the distinct elements problem. In *Proceedings of the Twenty-Ninth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS)*, pages 41–52, 2010. [20](#)
- [KP20] John Kallaugher and Eric Price. Separations and equivalences between turnstile streaming and linear sketching. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 1223–1236, 2020. [43](#)
- [KS92] Bala Kalyanasundaram and Georg Schnitger. The probabilistic communication complexity of set intersection. *SIAM J. Discret. Math.*, 5(4):545–557, 1992. [38](#)
- [KW11] Felix Krahmer and Rachel Ward. New and improved JohnsonLindenstrauss embeddings via the restricted isometry property. *SIAM Journal on Mathematical Analysis*, 43(3):1269–1281, 2011. [66](#), [70](#)
- [LN16] Kasper Green Larsen and Jelani Nelson. The johnson-lindenstrauss lemma is optimal for linear dimensionality reduction. In *Proceedings of the 43rd International Colloquium on Automata, Languages, and Programming (ICALP)*, pages 82:1–82:11, 2016. [61](#)

- [LN17] Kasper Green Larsen and Jelani Nelson. Optimality of the johnson-lindenstrauss lemma. In *Proceedings of the 58th IEEE Annual Symposium on Foundations of Computer Science (FOCS)*, pages 633–638, 2017. [61](#), [62](#)
- [LNNT16] Kasper Green Larsen, Jelani Nelson, Huy L. Nguyễn, and Mikkel Thorup. Heavy hitters via cluster-preserving clustering. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, 2016. [47](#)
- [LNW14] Yi Li, Huy L. Nguyen, and David P. Woodruff. Turnstile streaming algorithms might as well be linear sketches. In *Proceedings of the 46th Annual ACM Symposium on Theory of Computing (STOC)*, pages 174–183, 2014. [43](#)
- [LS17] Yin Tat Lee and He Sun. An SDP-based algorithm for linear-sized spectral sparsification. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 678–687, 2017. [82](#)
- [LS18] Yin Tat Lee and He Sun. Constructing linear-sized spectral sparsification in almost-linear time. *SIAM J. Comput.*, 47(6):2315–2336, 2018. [82](#)
- [Lum18] Jérémie O. Lumbroso. The story of HyperLogLog: How Flajolet processed streams with coin flips. *CoRR*, abs/1805.00612v2, 2018. [20](#)
- [McG14] Andrew McGregor. Graph stream algorithms: a survey. *SIGMOD Rec.*, 43(1):9–20, 2014. [49](#)
- [MM13] Xiangrui Meng and Michael W. Mahoney. Low-distortion subspace embeddings in input-sparsity time and applications to robust linear regression. In *Proceedings of the 45th ACM Symposium on Theory of Computing (STOC)*, pages 91–100, 2013. [82](#), [83](#)
- [MMR19] Konstantin Makarychev, Yury Makarychev, and Ilya P. Razenshteyn. Performance of johnson-lindenstrauss transform for  $k$ -means and  $k$ -medians clustering. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 1027–1038, 2019. [90](#), [93](#)
- [Mor78] Robert Morris. Counting large numbers of events in small registers. *Commun. ACM*, 21(10):840–842, 1978. [13](#), [15](#)
- [MPTJ07] Shahar Mendelson, Alain Pajor, and Nicole Tomczak-Jaegermann. Reconstruction and subgaussian operators in asymptotic geometric analysis. *Geometric and Functional Analysis*, 17:1248–1282, 2007. [111](#)
- [MRL98] Gurmeet Singh Manku, Sridhar Rajagopalan, and Bruce G. Lindsay. Approximate medians and other quantiles in one pass and with limited memory. In *Proceedings ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 426–435, 1998. [24](#), [27](#)
- [Nar18] Shyam Narayanan. Deterministic  $O(1)$ -approximation algorithms to 1-center clustering with outliers. In *APPROX*, pages 21:1–21:19, 2018. [77](#)
- [Nar19] Shyam Narayanan. 3-wise independent random walks can be slightly unbounded. *CoRR*, abs/1807.04910v2, 2019. [114](#)

- [Nel11] Jelani Nelson. *Sketching and Streaming High-Dimensional Vectors*. PhD thesis, Massachusetts Institute of Technology, June 2011. [54](#)
- [Nis92] Noam Nisan. Pseudorandom generators for space-bounded computation. *Comb.*, 12(4):449–461, 1992. [56](#)
- [NN13a] Jelani Nelson and Huy L. Nguyễn. OSNAP: faster numerical linear algebra algorithms via sparser subspace embeddings. In *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 117–126, 2013. [82](#), [83](#)
- [NN13b] Jelani Nelson and Huy L. Nguyễn. Sparsity lower bounds for dimensionality-reducing maps. In *Proceedings of the 45th ACM Symposium on Theory of Computing (STOC)*, pages 101–110, 2013. [66](#)
- [NN14] Jelani Nelson and Huy L. Nguyễn. Lower bounds for oblivious subspace embeddings. In *ICALP*, pages 883–894, 2014. [83](#)
- [Nol10] John P. Nolan. *Stable Distributions — Models for Heavy Tailed Data*. Birkhauser, 2010. [54](#)
- [NT08] Deanna Needell and Joel A. Tropp. CoSAMP: Iterative signal recovery from incomplete and inaccurate samples. *Applied and Computational Harmonic Analysis*, 2008. [101](#)
- [NY19] Jelani Nelson and Huacheng Yu. Optimal lower bounds for distributed and streaming spanning forest computation. In *Proceedings of the 30th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1844–1860, 2019. [49](#), [53](#)
- [NY20] Jelani Nelson and Huacheng Yu. Optimal bounds for approximate counting, 2020. [15](#)
- [Raz92] Alexander A. Razborov. On the distributional complexity of disjointness. *Theor. Comput. Sci.*, 106(2):385–390, 1992. [38](#)
- [Rud99] Mark Rudelson. Random vectors in the isotropic position. *J. Functional Analysis*, 164(1):60–72, 1999. [10](#)
- [RV08] Mark Rudelson and Roman Vershynin. On sparse reconstruction from Fourier and Gaussian measurements. *Communications on Pure and Applied Mathematics*, 61(8):1025–1045, 2008. [71](#), [100](#)
- [RV13] Mark Rudelson and Roman Vershynin. Hanson-Wright inequality and sub-gaussian concentration. *arXiv*, abs/1306.2872, 2013. [9](#)
- [Sar06] Tamás Sarlós. Improved approximation algorithms for large matrices via random projections. In *47th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2006), 21-24 October 2006, Berkeley, California, USA, Proceedings*, pages 143–152, 2006. [75](#), [79](#), [82](#), [84](#), [87](#)
- [SBAS04] Nisheeth Shrivastava, Chiranjeev Buragohain, Divyakant Agrawal, and Subhash Suri. Medians and beyond: new aggregation techniques for sensor networks. In *Proceedings of the 2nd International Conference on Embedded Networked Sensor Systems (SenSys)*, pages 239–249, 2004. [24](#)

- [She12] Alexander A. Sherstov. The communication complexity of gap hamming distance. *Theory Comput.*, 8(1):197–208, 2012. [40](#)
- [SS11] Daniel A. Spielman and Nikhil Srivastava. Graph sparsification by effective resistances. *SIAM J. Comput.*, 40(6):1913–1926, 2011. [80](#), [82](#)
- [SSS95] Jeanette P. Schmidt, Alan Siegel, and Aravind Srinivasan. Chernoff-Hoeffding bounds for applications with limited independence. *SIAM J. Discret. Math.*, 8(2):223–250, 1995. [51](#)
- [Tal96] Michel Talagrand. Majorizing measures, the generic chaining. *Ann. Probab.*, 24:1049–1103, 1996. [109](#)
- [Tal14] Michel Talagrand. *Upper and lower bounds for stochastic processes: modern methods and classical problems*. Springer, 2014. [111](#)
- [Tro11] Joel A. Tropp. Improved analysis of the subsampled randomized Hadamard transform. *Adv. Adapt. Data Anal.*, 3(1–2):115–126, 2011. Special issue on Sparse Representation of Data and Images. [82](#)
- [TSJ08] D. Sivakumar T. S. Jayram, Ravi Kumar. The one-way communication complexity of hamming distance. *Theory of Computing*, 4(1):129–135, 2008. [20](#), [40](#)
- [TZ12] Mikkel Thorup and Yin Zhang. Tabulation-based 5-independent hashing with applications to linear probing and second moment estimation. *SIAM J. Comput.*, 41(2):293–331, 2012. [54](#), [67](#)
- [Vid12] Thomas Vidick. A concentration inequality for the overlap of a vector on a large set, with application to the communication complexity of the gap-hamming-distance problem. *Chic. J. Theor. Comput. Sci.*, 2012. [40](#)
- [WC79] Mark N. Wegman and Larry Carter. New classes and applications of hash functions. In *Proceedings of the 20th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 175–182, 1979. [19](#)
- [Woo04] David P. Woodruff. Optimal space lower bounds for all frequency moments. In *Proceedings of the Fifteenth Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 167–175, 2004. [20](#), [40](#)
- [Woo14] David P. Woodruff. Sketching as a tool for numerical linear algebra. *Foundations and Trends in Theoretical Computer Science*, 10(1-2):1–157, 2014. [75](#)
- [Yu21] Huacheng Yu. Tight distributed sketching lower bound for connectivity. In *Proceedings of the 32nd Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2021. [49](#)