# Dimensionality Reduction — Notes 2

Jelani Nelson

`minilek@seas.harvard.edu`

August 11, 2015

## 1 Optimality theorems for JL

Yesterday we saw for MJL that we could achieve target dimension $m = O(\varepsilon^{-2} \log N)$, and for DJL we could achieve $m = O(\varepsilon^{-2} \log(1/\delta))$. The following theorems tell us that not much improvement is possible for MJL, and for DJL we have the optimal bound.

**Theorem 1** ([Alo03]). *For any $N > 1$ and $\varepsilon < 1/2$, there exist $N+1$ points in $\mathbb{R}^N$ such that achieving the MJL guarantee with distortion $1 + \varepsilon$ requires*

$$m \gtrsim \min\{n, \varepsilon^{-2}(\log N)/\log(1/\varepsilon)\}.$$

The $\log(1/\varepsilon)$ loss in the lower bound can be removed if the map must be linear.

**Theorem 2** ([LN14]). *For any $N > 1$ and $\varepsilon < 1/2$, there exist $N^{O(1)}$ points in $\mathbb{R}^N$ such that achieving the MJL guarantee with distortion $1 + \varepsilon$ using a linear map* requires
$$m \gtrsim \min\{n, \varepsilon^{-2} \log N\}.$$

For DJL, the upper bound is optimal.

**Theorem 3** ([JW13, KMN11]). *For any $\varepsilon, \delta < 1/2$, any DJL distribution must have*
$$m \gtrsim \min\{n, \varepsilon^{-2} \log(1/\delta)\}.$$

# 2 Example application: deterministic $\ell_1$ point query and heavy hitters

Yesterday's notes gives an example application of JL to $k$-means clustering. Today we give another application.

In the $\ell_1$ point query problem a vector $x \in \mathbb{R}^n$ is updated in the turnstile streaming model. A query is an index $i \in [n]$, and the response to the query should be a value $\tilde{x}$ such that $|x_i - \tilde{x}_i| \le \varepsilon \|x\|_1$. We show an argument of [NNW14] that the JL lemma implies the existence of a fixed deterministic $\Pi \in \mathbb{R}^{m \times n}$ with $m \lesssim \varepsilon^{-2} \log n$ such that such a $\tilde{x}$ can be recovered from $\Pi x$.

**Definition 1.** *We say that a matrix $\Pi$ with columns $\Pi_1, \ldots, \Pi_n$ is $\varepsilon$-incoherent if (1) $\|\Pi_i\|_2 = 1$ for all $i$, and (2) for all $i \ne j$, $|\langle \Pi_i, \Pi_j \rangle| \le \varepsilon$.*

**Theorem 4.** *If $\Pi \in \mathbb{R}^{m \times n}$ is $\varepsilon$-incoherent, then there is a polynomial time recovery algorithm $\mathcal{A}_\Pi$ such that given any $y = \Pi x$, if we define $\tilde{x} = \mathcal{A}_\Pi(y)$ then $\|\tilde{x} - x\|_\infty \le \varepsilon \|x\|_1$.*

*Proof.* The recovery algorithm will be $\mathcal{A}_\Pi(y) = \Pi^T y = \Pi^T \Pi x$. Thus

$$\tilde{x}_i = e_i^T \Pi^T \Pi x = \sum_{j=1}^{n} \langle \Pi_i, \Pi_j \rangle \, x_i = x_i + \sum_{i \ne j} \langle \Pi_i, \Pi_j \rangle \, x_i = x_i \pm \varepsilon \|x\|_1.$$

$\square$

Now we show the existence of such $\Pi$ with small $m$.

**Lemma 1.** $\forall \, \varepsilon \in (0, 1/2)$, *there is $\varepsilon$-incoherent $\Pi$ with $m \lesssim \varepsilon^{-2} \log n$.*

*Proof.* Consider the set of vectors $\{0, e_1, \ldots, e_n\}$. By the JL lemma, there exists $\Pi'$ with $O(\varepsilon^{-2} \log n)$ rows, and having columns $\Pi_i'$ such that (1) $\|\Pi_i'\|_2 = \|\Pi' e_i\|_2 = 1 \pm \varepsilon/3$, and (2) $\|\Pi_i' - \Pi_j'\|_2 = \|\Pi' e_i - \Pi' e_j\|_2 = (1 \pm \varepsilon/3)\sqrt{2}$ for all $i \ne j$. Let $\Pi$ be the matrix whose $i$th column is $\Pi_i'/\|\Pi_i'\|_2$. Then $\|\Pi_i\|_2 = 1$ for all $i$, as desired. Furthermore

$$2(1 \pm \varepsilon)^2 = \|\Pi_i - \Pi_j\|_2^2 = \|\Pi_i\|_2^2 + \|\Pi_j\|_2^2 - 2 \langle \Pi_i, \Pi_j \rangle .$$

Note $\|\Pi_i\|_2^2$ and $\|\Pi_j\|_2^2$ are both $1 \pm O(\varepsilon)$, implying $|\langle \Pi_i, \Pi_j \rangle| = O(\varepsilon)$. The lemma follows by applying this argument with $\varepsilon$ scaled down by a constant.

$\square$

# 3   Faster JL

Typically we have some high-dimensional computational geometry problem, and we use JL to speed up our algorithm in two steps: (1) apply a JL map $\Pi$ to reduce the problem to low dimension $m$, then (2) solve the lower-dimensional problem. As $m$ is made smaller, typically (2) becomes faster. However, ideally we would also like step (1) to be as fast as possible. In this section, we investigate two approaches to speed up the computation of $\Pi x$.

One of the analyses will make use of the following Bernstein bound.

**Theorem 5** (Bernstein's inequality). *Let $X_1, \ldots, X_n$ be independent random variables that are each at most $K$ almost surely, and where*

$$\sum_{i=1}^{n} \mathbb{E}(X_i - \mathbb{E}\, X_i)^2 = \sigma^2.$$

*Then for all $p \geq 1$*

$$\|\sum_{i=1}^{n} X_i - \mathbb{E} \sum_i X_i\|_p \lesssim \sigma\sqrt{p} + Kp.$$

*Proof.* Let $r_1, \ldots, r_n$ be independent Rademachers. Then

$$\|\sum_i (X_i - \mathbb{E}\, X_i)\|_p \leq 2 \cdot \|\sum_i r_i X_i\|_p \text{ (symmetrization)}$$

$$\lesssim \sqrt{p} \cdot \|(\sum_i X_i^2)^{1/2}\|_p \text{ (Khintchine)} \tag{1}$$

$$= \sqrt{p} \cdot \|\sum_i X_i^2\|_{p/2}^{1/2}$$

$$\leq \sqrt{p} \cdot \|\sum_i X_i^2\|_p^{1/2}$$

$$\leq \sigma\sqrt{p} + \sqrt{p} \cdot \|\sum_i X_i^2 - \mathbb{E} \sum_i X_i^2\|_p^{1/2} \text{ (triangle inequality)}$$

$$\leq \sigma\sqrt{p} + \sqrt{p} \cdot \|\sum_i r_i X_i^2\|_p^{1/2} \text{ (symmetrization)}$$

$$\leq \sigma\sqrt{p} + p^{3/4} \cdot \|(\sum_i X_i^4)^{1/2}\|_p^{1/2} \text{ (Khintchine)}$$

$$\leq \sigma\sqrt{p} + p^{3/4} \cdot \sqrt{K} \cdot \|(\sum_i X_i^2)^{1/2}\|_p^{1/2} \tag{2}$$

Defining $E = \|(\sum_i X_i^2)^{1/2}\|_p^{1/2}$ and comparing (1) with (2), for some constant $C > 0$

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

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

## 3.1   Sparse JL

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 $O(s \cdot \|x\|_0)$, where $\|x\|_0 = |\{i : x_i \neq 0\}|$. The goal is then to make $s, m$ as small as possible.

The following matrix $\Pi$ was introduced in [CCF04], and it was analyzed for DJL in [TZ12]. In this construction, one picks a hash function $h : [n] \to [m]$ from a pairwise independent family, and a function $\sigma : [n] \to \{-1, 1\}$ from a 4-wise independent family. Then for each $i \in [n]$, $\Pi_{h(i),i} = \sigma(i)$, 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 $s = 1$ 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. This was avoided yesterday 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, and if one unrolls the proof fully yesterday, one would find that yesterday's lecture obtained $\delta$ failure probability by using the Hanson-Wright $p$-norm bound for $p = \Theta(\log(1/\delta))$. That is to say, the improvement yesterday came from bounding a higher moment than $p = 2$ (i.e. Chebyshev).

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.

As previously, we assume $x \in \mathbb{R}^n$ has $\|x\|_2 = 1$. Our random $\Pi \in \mathbb{R}^{m \times n}$ 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.

4

- The $\eta_{r,i}$ are negatively correlated. That is, for any subset $S$ of $[m] \times [n]$, we have $\mathbb{E} \prod_{(r,i) \in S} \eta_{r,i} \leq \prod_{(r,i) \in S} \mathbb{E} \, \eta_{r,i} = (s/m)^{|S|}$.

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

**Theorem 6.** *As long as $m \simeq \varepsilon^{-2} \log(1/\delta)$ and $s \simeq \varepsilon m$,*

$$\forall x : \|x\|_2 = 1, \; \mathbb{P}_{\Pi}(|\|\Pi x\|_2^2 - 1| > \varepsilon) < \delta. \tag{3}$$

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

$$Z = \|\Pi x\|_2^2 - 1 = \frac{1}{s} \sum_{r=1}^{m} \sum_{i \neq j} \eta_{r,i} \eta_{r,j} \sigma_{r,i} \sigma_{r,j} x_i x_j \stackrel{\text{def}}{=} \sigma^T A_{x,\eta} \sigma,$$

Thus by Hanson-Wright

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

$A_{x,\eta}$ is a block diagonal matrix with $m$ blocks, where the $r$th block is $(1/s)x^{(r)}(x^{(r)})^T$ but with the diagonal zeroed out. Here $x^{(r)}$ is the vector with $(x^{(r)})_i = \eta_{r,i} x_i$. Now we just need to bound $\|\|A_{x,\eta}\|_F\|_p$ and $\|\|A_{x,\eta}\|\|_p$.

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

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

$$\|A_{x,\eta}\|_F^2 = \frac{1}{s^2} \sum_{i \neq j} x_i^2 x_j^2 \cdot Q_{i,j}.$$

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 $p$,

$$\begin{aligned}
\|\|A_{x,\eta}\|_F\|_p &= \|\|A_{x,\eta}\|_F^2\|_{p/2}^{1/2} \\
&\leq \|\frac{1}{s^2} \sum_{i \neq j} x_i^2 x_j^2 \cdot Q_{i,j}\|_p^{1/2} \\
&\leq \frac{1}{s} \left( \sum_{i \neq j} x_i^2 x_j^2 \cdot \|Q_{i,j}\|_p \right)^{1/2} \quad \text{(triangle inequality)}
\end{aligned}$$

5

$$\leq \frac{1}{\sqrt{m}}$$

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

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

We now show $\|Q_{i,j}\|_p \lesssim p$, for which we use Bernstein's inequality.

Suppose $\eta_{a_1,i}, \ldots, \eta_{a_s,i}$ are all 1, where $a_1 < a_2 < \ldots < a_s$. 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

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

$\square$

There are two natural distributions where $\eta$ satisfies the conditions for the SJLT. In the first, the columns are independent, and for each column $i$ $(\eta_{1,i}, \ldots, \eta_{m,i})$ is chosen uniformly at random from the $\binom{m}{s}$ vectors in $\{0,1\}^m$ having weight exactly $s$. A second distribution is the CountSketch of [CCF04]. In this distribution, we assume $s$ divides $m$, and the rows are partitioned arbitrarily into $s$ blocks each of equal size $m/s$ (e.g. the first $m/s$ rows, then the next $m/s$ rows, etc.). For each column $i$ and for each block $b$ with corresponding $\eta(b,i) = (\eta_{cm/s+1,i}, \ldots, \eta_{(c+1)m/s,i})$, we set $\eta(b,i) = e_j \in \mathbb{R}^{m/s}$ for a uniformly random $j \in [m/s]$. This is done independently across all $b, i$ pairs.

## 3.2 FFT-based approach

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 $O(n \log n)$, which is faster than the sparse JL approach when $x$ is sufficiently dense. Although we did not cover it this approach in lecture today, I am including a description here. They called their transformation the *Fast*

*Johnson-Lindenstrauss Transform (FJLT)*. A construction similar to theirs, which we will analyze here, is the $m \times n$ matrix $\Pi$ defined as

$$\Pi = \frac{1}{\sqrt{m}} SHD \qquad (4)$$

where $S$ is an $m \times n$ 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 $D = diag(\alpha)$ 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 $H^T H = I$ and $\max_{i,j} |H_{i,j}| \leq 1$. For example, $H$ can be the unnormalized Fourier matrix or Hadamard matrix. The original FJLT replaced $S$ with a random sparse matrix $P$, which has certain advantages; see Remark 1.

The motivation for the construction (4) is speed: $D$ can be applied in $O(n)$ time, $H$ in $O(n \log n)$ time (e.g. using the Fast Fourier Transform), and $S$ in $O(m)$ time. Thus, overall, applying $\Pi$ to any fixed vector $x$ takes $O(n \log n)$ time. Compare this with using a dense matrix of Rademachers, which takes $O(mn)$ time to apply.

We will show that for $m \gtrsim \varepsilon^{-2} \log(1/\delta) \log(1/(\varepsilon\delta))$, the random $\Pi$ described in (4) provides DJL. In fact we will analyze a slightly different construction in which $S$ is replaced by an $n \times n$ diagonal matrix $S_\eta$, $S_\eta = diag(\eta)$, where the entries of $\eta \in \{0, 1\}^n$ are independent with $\mathbb{E}\,\eta_i = 1/m$ (so $\Pi$ has $m$ rows in expectation). The proof to analyze the $\Pi$ from (4) is essentially identical. The proof we provide here is an adaptation of the proof of a more general theorem [CNW15, Theorem 9] to the current scenario.

**Theorem 7.** *Let $x \in \mathbb{R}^n$ be an arbitrary unit norm vector, and suppose $0 < \varepsilon, \delta < 1/2$. Also let $\Pi = S_\eta HD$ as described above with a number of rows equal to $m \gtrsim \varepsilon^{-2} \log(1/\delta) \log(1/(\varepsilon\delta))$. Then*

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

*Proof.* We use the moment method. Let $\eta'$ be an independent copy of $\eta$, and let $\sigma \in \{-1, 1\}^n$ be uniformly random. Write $z = HDx$ so that $\|\Pi x\|_2^2 = \sum_i \eta_i z_i^2$. Then

$$\|\frac{1}{m} \sum_{i=1}^{n} \eta_i z_i^2 - 1\|_p = \|\|\frac{1}{m} \sum_i \eta_i z_i^2 - 1\|_{L^p(\eta)}\|_{L^p(\alpha)} \qquad (5)$$

$$\leq \frac{2}{m} \cdot \|\| \sum_i \sigma_i \eta_i z_i^2 \|_{L^p(\eta)}\|_{L^p(\alpha)} \text{ (symmetrization)}$$

$$\leq \frac{2}{m} \cdot \| \sum_i \sigma_i \eta_i z_i^2 \|_p$$

$$\lesssim \frac{\sqrt{p}}{m} \cdot \|(\sum_i \eta_i z_i^4)^{1/2}\|_p \text{ (Khintchine)}$$

$$\leq \frac{\sqrt{p}}{m} \cdot \|(\max_i \eta_i |z_i|) \cdot (\sum_i \eta_i z_i^2)^{1/2}\|_p$$

$$\leq \frac{\sqrt{p}}{m} \cdot \| \max_i \eta_i z_i^2 \|_p^{1/2} \cdot \| \sum_i \eta_i z_i^2 \|_p^{1/2} \text{ (Cauchy-Schwarz)}$$

$$\leq \sqrt{\frac{p}{m}} \cdot \| \max_i \eta_i z_i^2 \|_p^{1/2} \cdot (\|\frac{1}{m} \sum_i \eta_i z_i^2 - 1\|_p^{1/2} + 1) \text{ (triangle inequality)}$$

$$(6)$$

We will now bound $\| \max_i \eta_i z_i^2 \|_p^{1/2}$. Define $q = \max\{p, \log m\}$ and note $\| \cdot \|_p \leq \| \cdot \|_q$. Then

$$\| \max_i \eta_i z_i^2 \|_q = \left( \mathbb{E}_{\alpha,\eta} \max_i \eta_i z_i^{2q} \right)^{1/q}$$

$$\leq \left( \mathbb{E}_{\alpha,\eta} \sum_i \eta_i z_i^{2q} \right)^{1/q}$$

$$= \left( \sum_i \mathbb{E}_{\alpha,\eta} \eta_i z_i^{2q} \right)^{1/q}$$

$$\leq \left( n \cdot \max_i \mathbb{E}_{\alpha,\eta} \eta_i z_i^{2q} \right)^{1/q}$$

$$= \left( n \cdot \max_i (\mathbb{E}_\eta \eta_i) \cdot (\mathbb{E}_\alpha z_i^{2q}) \right)^{1/q} \quad (\alpha, \eta \text{ independent})$$

$$= \left( m \cdot \max_i \mathbb{E}_\alpha z_i^{2q} \right)^{1/q}$$

$$\leq 2 \cdot \max_i \|z_i^2\|_q \ (m^{1/q} \leq 2 \text{ by choice of } q)$$

$$= 2 \cdot \max_i \|z_i\|_{2q}^2$$

$$\lesssim q \text{ (Khintchine)} \qquad\qquad (7)$$

8

Eq. (7) uses that $H$ is an unnormalized bounded orthonormal system.

Defining $E = \|\frac{1}{m} \sum_i \eta_i z_i^2 - 1\|_p^{1/2}$ and combining (5), (6), (7), we find that for some constant $C > 0$

$$E^2 - C\sqrt{\frac{pq}{m}} E - C\sqrt{\frac{pq}{m}} \leq 0,$$

implying $E^2 \lesssim \max\{\sqrt{pq/m}, pq/m\}$. By the Markov inequality

$$\mathbb{P}(|\|\Pi x\|_2^2 - 1| > \varepsilon) \leq \varepsilon^{-p} \cdot E^{2p},$$

and thus to achieve the theorem statement it suffices to set $p = \log(1/\delta)$ then choose $m \gtrsim \varepsilon^{-2} \log(1/\delta) \log(m/\delta)$. $\qquad\square$

**Remark 1.** 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 2.** The usual analysis for the FJLT, such as the approach in [AC09], would achieve a bound on $m$ of $O(\varepsilon^{-2} \log(1/\delta) \log(n/\delta))$. Such analyses operate by, using the notation of the proof of Theorem 7, first conditioning on $\|z\|_\infty \lesssim \sqrt{\log(n/\delta)}$ (which happens with probability at least $1 - \delta/2$ by the Khintchine inequality), then finishing the proof using Bernstein's inequality. In our proof above, we improved this dependence on $n$ to a dependence on the smaller quantity $m$ by avoiding any such conditioning.

# References

[AC09]   Nir Ailon and Bernard Chazelle. The fast Johnson–Lindenstrauss transform and approximate nearest neighbors. *SIAM J. Comput.*, 39(1):302–322, 2009.

[Alo03]   Noga Alon. Problems and results in extremal combinatorics–I. *Discrete Mathematics*, 273(1-3):31–53, 2003.

[CCF04]   Moses Charikar, Kevin C. Chen, and Martin Farach-Colton. Finding frequent items in data streams. *Theor. Comput. Sci.*, 312(1):3–15, 2004.

[CNW15]  Michael B. Cohen, Jelani Nelson, and David P. Woodruff. Optimal approximate matrix product in terms of stable rank. *CoRR*, abs/1507.02268, 2015.

[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.

[KMN11]  Daniel M. Kane, Raghu Meka, and Jelani Nelson. Almost optimal explicit Johnson-Lindenstrauss families. In *RANDOM*, pages 628–639, 2011.

[KN14]    Daniel M. Kane and Jelani Nelson. Sparser Johnson-Lindenstrauss transforms. *Journal of the ACM*, 61(1):4, 2014.

[LN14]    Kasper Green Larsen and Jelani Nelson. The Johnson-Lindenstrauss lemma is optimal for linear dimensionality reduction. *CoRR*, abs/1411.2404, 2014.

[NNW14]  Jelani Nelson, Huy L. Nguyễn, and David P. Woodruff. On deterministic sketching and streaming for sparse recovery and norm estimation. *Linear Algebra and its Applications, Special Issue on Sparse Approximate Solution of Linear Systems*, 441:152–167, 2014.

[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.