

Fast Dimension Reduction Using Rademacher Series on Dual BCH Codes

Nir Ailon*

Edo Liberty†

Abstract

The Fast Johnson-Lindenstrauss Transform (FJLT) was recently discovered by Ailon and Chazelle as a novel technique for performing fast dimension reduction with small distortion from ℓ_2^d to ℓ_2^k in time $O(\max\{d \log d, k^3\})$. For k in $[\Omega(\log d), O(d^{1/2})]$ this beats time $O(dk)$ achieved by naive multiplication by random dense matrices, an approach followed by several authors as a variant of the seminal result by Johnson and Lindenstrauss (JL) from the mid 80's. In this work we show how to significantly improve the running time to $O(d \log k)$ for $k = O(d^{1/2-\delta})$, for any arbitrary small fixed δ . This beats the better of FJLT and JL. Our analysis uses a powerful measure concentration bound due to Talagrand applied to Rademacher series in Banach spaces (sums of vectors in Banach spaces with random signs). The set of vectors used is a real embedding of dual BCH code vectors over $GF(2)$. We also discuss the number of random bits used and reduction to ℓ_1 space.

The connection between geometry and discrete coding theory discussed here is interesting in its own right and may be useful in other algorithmic applications as well.

1 Introduction

Applying random matrices is by now a well known technique for reducing dimensionality of vectors in Euclidean space while preserving certain properties (most notably distance information). Beginning with the classic work of Johnson and Lindenstrauss [16], who used projections onto random subspaces, other variants of the technique using different distributions are known [1, 9, 15, 22] and have been used in many algorithms [18, 20, 3, 13, 26, 24, 11].

In all the variants of this idea, a fixed unit length vector $x \in \mathbf{R}^d$ is mapped onto \mathbf{R}^k ($k < d$) via a random linear mapping Φ from a carefully chosen distribution. A measure concentration principle is used to show that

the distribution of the norm estimator error $|\|\Phi x\|_2 - 1|$ in a small neighborhood of 0 is dominated by a Gaussian of standard deviation $\Omega(k^{-1/2})$, uniformly for all x and independent of d . The distribution of Φ need not even be rotationally invariant. When used in an algorithm, k is often chosen as $O(\varepsilon^{-2} \log n)$ so that a union bound ensures that the error is smaller than a fixed ε simultaneously for all n vectors in some fixed input set. Noga Alon proved [2] that this choice of k is essentially optimal and cannot be significantly reduced.

It makes sense to abstract the definition of a distribution of mappings that can be used for dimension reduction in the above sense. We will say that such a mapping has the Johnson-Lindenstrauss property (JLP), named after the authors of the first such construction (we make an exact definition of this property in Section 2). In view of Ailon and Chazelle's FJLT result [1], it is natural to ask about the computational complexity of applying a mapping drawn from a JLP distribution on a vector. The resources considered here are time and randomness. Ailon et al showed that reduction from d dimensions to k dimensions can be performed in time $O(\max\{d \log d, k^3\})$, beating the naïve $O(kd)$ time implementation of JL for k in $\omega(\log d)$ and $o(d^{1/2})$. Similar bounds were found in [1] for reducing onto ℓ_1 (Manhattan) space, but with quadratic (not cubic) dependence on k . From recent work by Matousek [22] it can be shown, by replacing gaussian distributions with ± 1 's, that Ailon and Chazelle's algorithm for the Euclidean case requires $O(\max\{d, k^3\})$ random bits in the Euclidean case.

1.1 Our Results This work contains several contributions. We summarize them for the Euclidean case in Table 1.1 for convenience. The first (in Section 7) is a simple trick that can be used to reduce the running time of FJLT [1] to $O(\max\{d \log k, k^3\})$, hence making it better than the naïve algorithm for small k (first column in the table). In typical applications, the running time translates to $O(d \log \log n)$, where n is the number of points we simultaneously want to reduce (assuming $n = 2^{O(d^{1/3})}$).

The main contribution (Sections 5-6) is improving

*Institute for Advanced Study, Princeton NJ and Google Research, New-York NY (nailon@google.com)

†Yale University, New Haven CT (edo.liberty@yale.edu)
Supported by AFOSR, and NGA.

	Fast — — — — — → Slow		
k in $o(\log d)$	This work	JL	FJLT
k in $\omega(\log d)$ and $o(\text{poly}(d))$	This work	FJLT	JL
k in $\Omega(\text{poly}(d))$ and $o((d \log d)^{1/3})$	This work, FJLT		JL
k in $\omega((d \log d)^{1/3})$ and $O(d^{1/2-\delta})$	This work	FJLT	JL

Table 1.1: Schematic comparison of asymptotic running time of this work, Ailon and Chazelle’s work [1] (FJLT) and a naïve implementation of Johnson-Lindenstrauss (JL), or variants thereof.

the case of "large k " (rightmost column in the Table 1.1). We use tools from the theory of probability and norm interpolation in Banach spaces (Section 3) as well as the theory of error correcting codes (Section 4) to construct a distribution on matrices satisfying JLP that can be applied in time $O(d \log d)$ (note that, in this case, $\log d = O(\log k)$). The ideas used in our constructions take advantage of advanced ideas from different classical theories. These ideas provide a new algorithmic application of error correcting codes, an extremely useful tool in theoretical computer science with applications in both complexity and algorithms (a good overview can be found in [25]; a more recent example in [17]).

A note on "large k ": As stated above, k is typically $O(\varepsilon^{-2} \log n)$, where ε is a desired distortion bound and n is the number of vectors we seek to reduce. Although $\log n$ is typically small (logarithmic in input size), in various applications, especially in scientific computation, ε^{-2} may be large. This case is therefore important to study.

It is illustrative to point out an apparent weakness in [1] that was a starting point of our work. The main tool used there was to multiply the input vector x by a random sign change matrix followed by a Fourier transform, resulting in a vector y . It is claimed that $\|y\|_\infty$ is small (in other words, the "information" is spread out evenly among the coordinates). By a convexity argument the "worst case" y (assuming only the ℓ_∞ bound) is a *uniformly supported* vector, namely,

a vector in which the absolute value of the coordinates in its (small) support are all equal. Intuitively, such a vector is extremely unlikely. In this work we consider other norms.

It is likely that the techniques we develop here can be used in conjunction with very recent research on explicit embeddings of ℓ_2 in ℓ_1 [23, 14, 4] as well as research on fast approximate linear algebraic scientific computation [24, 10, 6, 7, 8, 27].

2 Preliminaries

We use ℓ_p^d to denote d dimensional real space equipped with the norm $\|x\| = \|x\|_p = \left(\sum_{i=1}^d |x_i|^p\right)^{1/p}$, where $1 \leq p < \infty$ and $\|x\|_\infty = \max\{|x_i|\}$. The dual norm index q is defined by the solution to $1/q + 1/p = 1$. We remind the reader that $\|x\|_p = \sup_{\|y\|_q=1} x^T y$. For a real $k \times d$ matrix A , the matrix norm $\|A\|_{p_1 \rightarrow p}$ is defined as the operator norm of $A : \ell_{p_1}^d \rightarrow \ell_p^k$ or:

$$\|A\|_{p_1 \rightarrow p} = \sup_{\substack{x \in \ell_{p_1}^d \\ \|x\|=1}} \|Ax\|_p = \sup_{\substack{y \in \ell_q^k \\ \|y\|=1}} \sup_{\substack{x \in \ell_{p_1}^d \\ \|x\|=1}} y^T Ax.$$

In what follows we use d to denote the original dimension and $k < d$ the target (reduced) dimension. The input vector will be $x = (x_1, \dots, x_d)^T \in \ell_2^d$. Since we only consider linear reductions we will assume without loss of generality that $\|x\|_2 = 1$.

DEFINITION 2.1. *A distribution $\mathcal{D}(d, k)$ on $k \times d$ real matrices ($k \leq d$) has the Johnson-Lindenstrauss property (JLP) with respect to a norm index p , if for any unit vector $x \in \ell_2^d$ and $0 \leq \varepsilon < 1/2$,*

$$(2.1) \quad \Pr_{A \sim \mathcal{D}_{d,k}} [\|Ax\|_p - 1| > \varepsilon] \leq c_1 e^{-c_2 k \varepsilon^2}$$

for some global $c_1, c_2 > 0$.

(A similar definition was given in [24].) In this work, we study the cases $p = 1$ (*Manhattan JLP*) and $p = 2$ (*Euclidean JLP*). We make a few technical remarks about Definition 2.1:

- For most dimension reduction applications $k = \Omega(\varepsilon^{-2})$, so the constant c_1 can be "swallowed" by c_2 , but we prefer to keep it here to avoid writing $O(e^{-\Omega(k\varepsilon^2)})$ and for generality.
- The definition is robust with respect to bias of $O(k^{-1/2})$. More precisely, if we prove $\Pr[\mu - \varepsilon \leq \|Ax\|_p \leq \mu + \varepsilon] \geq 1 - c_1 e^{-c_2 k \varepsilon^2}$ for some μ satisfying $|\mu - 1| = O(k^{-1/2})$, then this would imply (2.1), with possibly different constants. We will use this observation in what follows.

Recall that a Walsh-Hadamard matrix H_d is a $d \times d$ orthogonal matrix with $H_d(i, j) = 2^{-d/2}(-1)^{\langle i, j \rangle}$ for all $i, j \in [0, d-1]$, where $\langle i, j \rangle$ is dot product (over \mathbb{F}_2) of i, j viewed as $(\log d)$ -bit vectors. The matrix encodes the Fourier transform over the binary hypercube. It is well known that $x \mapsto H_d x \in \ell_2^d$ can be computed in time $O(d \log d)$ for any $x \in \ell_2^d$, and that the mapping is isomorphic.

DEFINITION 2.2. *A matrix $A \in \mathbf{R}^{m \times d}$ is a code matrix if every row of A is equal to some row of H multiplied by $\sqrt{d/m}$.*

The normalization is chosen so that columns have Euclidean norm 1.

2.1 Statement of our Theorems The main contribution is in Theorem 2.2 below.

THEOREM 2.1. *For any code matrix A of size $k \times d$ for $k < d$, the mapping $x \mapsto Ax$ can be computed in time $O(d \log k)$.*

Clearly this theorem is interesting only for $\log k = o(\log d)$, because otherwise the Walsh-Hadamard transform followed by projection onto a subset of the coordinates can do this in time $O(d \log d)$, by definition of a code matrix. As a simple corollary, the running time of the algorithms in [1] can be reduced to $O(\max\{d \log k, k^3\})$, because effectively what they do is multiply the input x (after a random sign change) by a code matrix of size $O(k^3) \times d$ and then manipulate the outcome in time $O(k^3)$. This gives the left column of Table 1.1. We omit the details of this result and refer the reader to [1, 22].

THEOREM 2.2. *Let $\delta > 0$ be some arbitrarily small constant. For any d, k satisfying $k \leq d^{1/2-\delta}$ there exists an algorithm constructing a random matrix A of size $k \times d$ satisfying JLP, such that the time to compute $x \mapsto Ax$ for any $x \in \mathbf{R}^d$ is $O(d \log k)$. The construction uses $O(d)$ random bits and applies to both the Euclidean and the Manhattan cases.*

We will prove a slightly weaker running time of $O(d \log d)$ below, and provide a sketch for reducing it to $O(d \log k)$, where the full details of the improvement are deferred to Appendix A. This improvement is interesting for small k , and provides a unified solution for all $k \leq d^{1/2-\delta}$, though the small k case can also be taken care of using Theorem 2.1 above in conjunction with FJLT [1]. The main contribution of theorem 2.1, of course, is in getting rid of the term k^3 in the running time of FJLT.

3 Tools from Banach Spaces

The following is known as an interpolation theorem in the theory of Banach spaces. For a proof, refer to [5].

THEOREM 3.1. [Riesz-Thorin] *Let A be an $m \times d$ real matrix, and assume $\|A\|_{p_1 \rightarrow r_1} \leq C_1$ and $\|A\|_{p_2 \rightarrow r_2} \leq C_2$ for some norm indices p_1, r_1, p_2, r_2 . Let λ be a real number in the interval $[0, 1]$, and let p, r be such that $1/p = \lambda(1/p_1) + (1-\lambda)(1/p_2)$ and $1/r = \lambda(1/r_1) + (1-\lambda)(1/r_2)$. Then $\|A\|_{p \rightarrow r} \leq C_1^\lambda C_2^{1-\lambda}$.*

THEOREM 3.2. [Hausdorff-Young] *For norm index $1 \leq p \leq 2$, $\|H\|_{p \rightarrow q} \leq d^{-1/p+1/2}$, where q is the dual norm index of p .*

(This theorem is usually stated with respect to the Fourier operator for functions on the real line or on the circle, and is a simple application of Riesz-Thorin by noticing that $\|H\|_{2 \rightarrow 2} = 1$ and $\|H\|_{1 \rightarrow \infty} = d^{-1/2}$.)

Let M be a real $m \times d$ matrix, and let $z \in \mathbf{R}^d$ be a random vector with each z_i distributed uniformly and independently over $\{\pm 1\}$. The random vector $Mz \in \ell_p^m$ is known as a *Rademacher* random variable. A nice exposition of concentration bounds for Rademacher variables is provided in Chapter 4.7 of [19] for more general Banach spaces. For our purposes, it suffices to review the result for finite dimensional ℓ_p space. Consider the norm $Z = \|Mz\|_p$ (we say that " Z is the norm of a Rademacher random variable in ℓ_p^d corresponding to M "). We associate two numbers with Z ,

- The deviation σ , defined as $\|M\|_{2 \rightarrow p}$, and
- The median μ of Z .

THEOREM 3.3. *For any $t \geq 0$, $\Pr[|Z - \mu| > t] \leq 4e^{-t^2/(8\sigma^2)}$.*

The theorem is a simple consequence of a powerful theorem of Talagrand (Chapter 1, [19]) on measure concentration of functions on $\{-1, +1\}^d$ extendable to convex functions on ℓ_2^d with bounded Lipschitz norm.

4 Tools from Error Correcting Codes

Let A be a code matrix, as defined above. The columns of A can be viewed as vectors over \mathbb{F}_2 under the usual transformation ($(+) \rightarrow 0, (-) \rightarrow 1$). Clearly, the set of vectors thus obtained are closed under addition, and hence constitute a linear subspace of \mathbb{F}_2^m . Conversely, any linear subspace V of \mathbb{F}_2^m of dimension ν can be encoded as an $m \times 2^\nu$ code matrix (by choosing some ordered basis of V). We will borrow well known constructions of subspaces from coding theory, hence the terminology. Incidentally, note that H_d encodes

the Hadamard code, equivalent to a dual BCH code of designed distance 3.

DEFINITION 4.1. A code matrix A of size $m \times d$ is a -wise independent if for each $1 \leq i_1 < i_2 < \dots < i_a \leq m$ and $(b_1, b_2, \dots, b_a) \in \{+1, -1\}^a$, the number of columns $A^{(j)}$ for which $(A_{i_1}^{(j)}, A_{i_2}^{(j)}, \dots, A_{i_a}^{(j)}) = (b_1, b_2, \dots, b_a)$ is exactly $d/2^a$.

LEMMA 4.1. There exists a 4-wise independent code matrix of size $k \times f_{\text{BCH}}(k)$, where $f_{\text{BCH}}(k) = \Theta(k^2)$.

The family of matrices is known as binary dual BCH codes of designed distance 5. Details of the construction can be found in [21].

5 Reducing to Euclidean Space for $k \leq d^{1/2-\delta}$

Assume $\delta > 0$ is some arbitrarily small constant. Let B be a $k \times d$ matrix with Euclidean unit length columns, and D a random $\{\pm 1\}$ diagonal matrix. Let $Y = \|BDx\|_2$. Our goal is to get a concentration bound of Y around 1. Notice that $E[Y^2] = 1$. In order to use Theorem 3.3, we let M denote the $k \times d$ matrix with its i 'th column $M^{(i)}$ being $x_i B^{(i)}$, where $B^{(i)}$ denotes the i 'th column of B . Clearly Y is the norm of a Rademacher random variable in ℓ_2^k corresponding to M . We estimate the deviation σ and median μ , as defined in Section 3.

$$\begin{aligned} \sigma &= \|M\|_{2 \rightarrow 2} = \sup_{\substack{y \in \ell_2^k \\ \|y\|=1}} \|y^T M\|_2 \\ &= \sup \left(\sum_{i=1}^d x_i^2 (y^T B^{(i)})^2 \right)^{1/2} \\ (5.2) \quad &\leq \|x\|_4 \sup \left(\sum_{i=1}^d (y^T B^{(i)})^4 \right)^{1/4} \\ &= \|x\|_4 \|B^T\|_{2 \rightarrow 4}. \end{aligned}$$

(The inequality is Cauchy-Schwartz.) To estimate the median, μ , we compute

$$\begin{aligned} E[(Z - \mu)^2] &= \int_0^\infty \Pr[(Z - \mu)^2 > s] ds \\ &\leq \int_0^\infty 4e^{-s/(8\sigma^2)} ds = 32\sigma^2. \end{aligned}$$

The inequality is an application of Theorem 3.3. Recall that $E[Z^2] = 1$. Also, $E[Z] = E[\sqrt{Z^2}] \leq \sqrt{E[Z^2]} = 1$ (by Jensen). Hence $E[(Z - \mu)^2] = E[Z^2] - 2\mu E[Z] + \mu^2 \geq 1 - 2\mu + \mu^2 = (1 - \mu)^2$. Combining, $|1 - \mu| \leq \sqrt{32}\sigma$. We conclude,

COROLLARY 5.1. For any $t \geq 0$,

$$\Pr[|Z - 1| > t] \leq c_3 \exp\{-c_4 t^2 / (\|x\|_4^2 \|B^T\|_{2 \rightarrow 4}^2)\},$$

for some global $c_3, c_4 > 0$.

In order for the distribution of BD to satisfy JLP, we need to have $\sigma = O(k^{-1/2})$. This requires controlling both $\|B^T\|_{2 \rightarrow 4}$ and $\|x\|_4$. We first show how to design a matrix B that is both efficiently computable and has a small norm. The latter quantity is adversarial and cannot be directly controlled, but we are allowed to manipulate x by applying a (random) orthogonal matrix Φ without losing any information.

5.1 Bounding $\|B^T\|_{2 \rightarrow 4}$ Using BCH Codes

LEMMA 5.1. Assume B is a $k \times d$ 4-wise independent code matrix. Then $\|B^T\|_{2 \rightarrow 4} \leq (3d)^{1/4} k^{-1/2}$.

Proof. For $y \in \ell_2^k$, $\|y\| = 1$,

$$\begin{aligned} (5.3) \quad \|y^T B\|_4^4 &= d E_{j \in [d]} [(y^T B^{(j)})^4] \\ &= dk^{-2} \sum_{i_1, i_2, i_3, i_4=1}^k E_{b_1, b_2, b_3, b_4} [y_{i_1} y_{i_2} y_{i_3} y_{i_4} b_1 b_2 b_3 b_4] \\ &= dk^{-2} (3\|y\|_2^4 - 2\|y\|_4^4) \leq 3dk^{-2}, \end{aligned}$$

where b_1, b_2, b_3, b_4 are random $\{+1, -1\}$ variables. We now use the BCH codes. Let B_k denote the $k \times f_{\text{BCH}}(k)$ matrix from the Lemma 4.1 (we assume here that $k = 2^a - 1$ for some integer a ; This is harmless because otherwise we can reduce onto some $k' = 2^a - 1$ such that $k/2 \leq k' \leq k$ and pad the output with $k - k'$ zeros). In order to construct a matrix B of size $k \times d$ for $k \leq d^{1/2-\delta}$, we first make sure that d is divisible by $f_{\text{BCH}}(k)$ (by at most multiplying d by a constant factor and padding with zeros), and then define B to be $d/f_{\text{BCH}}(k)$ copies of B_k side by side. Clearly B remains 4-wise independent. Note that B may no longer be a code matrix, but $x \mapsto Bx$ is computable in time $O(d \log k)$ by performing $d/f_{\text{BCH}}(k)$ Walsh transforms on blocks of size $f_{\text{BCH}}(k)$.

5.2 Controlling $\|x\|_4$ for $k < d^{1/2-\delta}$ We define a randomized orthogonal transformation Φ that is computable in $O(d \log d)$ time and succeeds with probability $1 - O(e^{-k})$ for all $k < d^{1/2-\delta}$. Success means that $\|\Phi x\|_4 = O(d^{-1/4})$. (Note: Both big- O 's hide factors depending on δ). Note that this construction gives a running time of $O(d \log d)$. We discuss later how to do this for arbitrarily small k with running time $O(d \log k)$.

The basic building block is the product HD' , where $H = H_d$ is the Walsh-Hadamard matrix and D' is a

diagonal matrix with random i.i.d. uniform $\{\pm 1\}$ on the diagonal. Note that this random transformation was the main ingredient in [1]. Let $H^{(i)}$ denote the i 'th column of H .

We are interested in the random variable $X = \|HD'x\|_4$. We define M as the $d \times d$ matrix with the i 'th column $M^{(i)}$ being $x_i H^{(i)}$, we let $p = 4$ ($q = 4/3$), and notice that X is the norm of the Rademacher random variable in ℓ_4^d corresponding to M (using the notation of Section 3). We compute the deviation σ ,

$$\begin{aligned}
\sigma &= \|M\|_{2 \rightarrow 4} = \|M^T\|_{4/3 \rightarrow 2} \\
&= \sup_{\substack{y \in \ell_{4/3}^k \\ \|y\|_{4/3} = 1}} \left(\sum_i x_i^2 (y^T H^{(i)})^2 \right)^{1/2} \\
&\leq \left(\sum_i x_i^4 \right)^{1/4} \sup \left(\sum_i (y^T H^{(i)})^4 \right)^{1/4} \\
&= \|x\|_4 \|H^T\|_{\frac{4}{3} \rightarrow 4}.
\end{aligned} \tag{5.4}$$

(Note that $H^T = H$.) By the Hausdorff-Young theorem, $\|H\|_{\frac{4}{3} \rightarrow 4} \leq d^{-1/4}$. Hence, $\sigma \leq \|x\|_4 d^{-1/4}$. We now get by Theorem 3.3 that for all $t \geq 0$,

$$(5.5) \quad \Pr[|\|HD'x\|_4 - \mu| > t] \leq 4e^{-t^2/(8\|x\|_4^2 d^{-1/2})},$$

where μ is a median of X .

CLAIM 5.1. $\mu = O(d^{-1/4})$.

Proof. To see the claim, notice that, for each separate coordinate, $E[(HD'x)_i^4] = O(d^{-2})$ and then use linearity of expectation to get $E[\|HD'x\|_4^4] = O(d^{-1})$. By Jensen's inequality, $E[\|HD'x\|_4^b] \leq E[\|HD'x\|_4^4]^{b/4} = O(d^{-b/4})$ for $b = 1, 2, 3$. Now

$$\begin{aligned}
E[(\|HD'x\|_4 - \mu)^4] &= \int_0^\infty \Pr[(\|HD'x\|_4 - \mu)^4 > s] ds \\
&\leq \int_0^\infty 4e^{-s^{1/2}/(8\|x\|_4^2 d^{-1/2})} ds \\
&= O(d^{-1}).
\end{aligned}$$

This implies that $\gamma_1 d^{-1} - \gamma_2 d^{-3/4} \mu + \gamma_3 d^{-2/4} \mu^2 - \gamma_4 d^{-1/4} \mu^3 + \mu^4 \leq \gamma_5 d^{-1}$, where γ_i is a global constant for $i = 1, 2, 3, 4, 5$. Clearly this implies the statement of the claim.

Let c_9 be such that $\mu_4 \leq c_9 d^{-1/4}$. We weaken inequality (5.5) using the last claim to obtain the following convenient form:

$$(5.6) \quad \Pr[\|HD'x\|_4 > c_9 d^{-1/4} + t] \leq 4e^{-t^2/(8\|x\|_4^2 d^{-1/2})}.$$

In order to get a desired failure probability of $O(e^{-k})$ set $t = c_8 k^{1/2} \|x\|_4 d^{-1/4}$. For $k < d^{1/2-\delta}$ this

gives $t < c_8 d^{-\delta/2} \|x\|_4$. In other words, with probability $1 - O(e^{-k})$ we get

$$\|HD'x\|_4 \leq c_9 d^{-1/4} + c_8 d^{-\delta/2} \|x\|_4.$$

Now compose this r times: Take independent random diagonal $\{\pm 1\}$ matrices $D' = D^{(1)}, D^{(2)}, \dots, D^{(r)}$ and define $\Phi_d^{(r)} = HD^{(r)} HD^{(r-1)} \dots HD^{(1)}$. Using a union bound on the conditional failure probabilities, we easily get:

LEMMA 5.2. [ℓ_4 **reduction for** $k < d^{1/2-\delta}$] With probability $1 - O(e^{-k})$

$$(5.7) \quad \|\Phi^{(r)}x\|_4 = O(d^{-1/4})$$

for $r = \lceil 1/2\delta \rceil$.

Note that the constant hiding in the bound (5.7) is exponential in $1/\delta$.

Combining the above, the random transformation $A = BD\Phi^{(r)}$ has Euclidean JLP for $k < d^{1/2-\delta}$, and can be applied to a vector in time $O(d \log d)$. This proves the Euclidean case of Theorem 2.2.

5.3 Reducing the Running Time to $O(d \log k)$

We now explain how to reduce the running time to $O(d \log k)$, using the new tools developed here. This provides a unified solution to the problem of designing efficient Johnson-Lindenstrauss projections for all k up to $d^{1/2-\delta}$. Recall that in the construction of B we placed $d/f_{\text{BCH}}(k)$ copies of the same code matrix B_k of size $k \times f_{\text{BCH}}(k)$ side by side. It turns out that we can apply this "decomposition" of coordinates to $\Phi^{(r)}$. Indeed, let $I_j \subseteq [d]$ denote the j 'th block of $\beta = f_{\text{BCH}}(k)k^\delta$ consecutive coordinates (assume that β is an integer that divides d). For a vector $y \in \ell_p^d$, let $y_{I_j} \in \ell_p^\beta$ denote the projection of y onto the set of coordinates I_j . Now, instead of using $\Phi^{(r)} = \Phi_d^{(r)}$ as above, we use a block-diagonal $d \times d$ matrix comprised of d/β $\beta \times \beta$ blocks each drawn from the same distribution as $\Phi_\beta^{(r)}$. Clearly the running time of the block-diagonal matrix is $O(d \log k)$, by applying the Walsh transform independently on each block (recall that $\beta = f_{\text{BCH}}(k)k^\delta = O(k^{2+\delta})$).

In order to see why this still works, one needs to repeat the above proofs using a family of norms $\|\cdot\|_{(p_1, p_2)}$ indexed by two norm indices p_1, p_2 and defined as $\|x\|_{(p_1, p_2)} = \left(\sum_{j=1}^{d/\beta} \|x_{I_j}\|_{p_1}^{p_2} \right)^{1/p_2}$. We discuss this in detail in Appendix A.

6 Reducing to Manhattan Space for $k < d^{1/2-\delta}$

We sketch this simpler case. As we did for the Euclidean case, we start by studying the random variable $W \in \ell_1^k$

defined as $W = \|k^{1/2}BDx\|_1$ for B as described in Section 5 and D a random ± 1 -diagonal matrix. In order to characterize the concentration of W (the norm of a Rademacher r.v. in ℓ_1^k) we compute the deviation σ , and estimate a median μ . As before, we set M to be the $k \times d$ matrix with the i 'th column being $k^{1/2}B^{(i)}x_i$.

(6.8)

$$\begin{aligned} \sigma &= \sup_{\substack{y \in \ell_k^\infty \\ \|y\|=1}} \|y^T M\|_2 = \sup \left(k \sum_{i=1}^d x_i^2 (y^T B^{(i)})^2 \right)^{1/2} \\ &\leq \sup k^{1/2} \|x\|_4 \|y^T B^{(i)}\|_4 = k^{1/2} \|x\|_4 \|B^T\|_{\infty \rightarrow 4} \end{aligned}$$

Using the tools developed in the Euclidean case, we can reduce $\|x\|_4$ to $O(d^{-1/4})$ with probability $1 - O(e^{-k})$ using $\Phi_r(d)$, in time $O(d \log d)$ (in fact, $O(d \log k)$ using the improvement from Appendix A). Also we already know from Section 5.1 that $\|B^T\|_{2 \rightarrow 4} = O(d^{1/4}k^{-1/2})$ if B is comprised of $k \times f_{\text{BCH}}(k)$ dual BCH codes (of designed distance 5) matrices side by side (assume $f_{\text{BCH}}(k)$ divides d). Since $\|y\|_\infty \geq k^{-1/2}\|y\|_2$ for any $y \in \ell_k$, we conclude that $\|B^T\|_{\infty \rightarrow 4} = O(d^{1/4})$. Combining, we get $\sigma = O(k^{1/2})$. We now estimate the median μ of W .

In order to calculate μ we first calculate $E(W) = kE[|P|]$ where P is any single coordinate of $k^{1/2}BDx$. We follow (almost exactly) a proof by Matousek in [22] where he uses a quantitative version of the Central Limit Theorem by König, Schütt, and Tomczak [12].

LEMMA 6.1. [König-Schütt-Tomczak] *Let $z_1 \dots z_d$ be independent symmetric random variables with $\sum_{i=1}^d E[z_i^2] = 1$, let $F(t) = \Pr[\sum_{i=1}^d z_i < t]$, and let $\bar{\varphi}(t) = \frac{1}{2\pi} \int_{-\infty}^t e^{-x^2/2} dx$. Then*

$$|F(t) - \bar{\varphi}(t)| \leq \frac{C}{1 + |t|^3} \sum_{i=1}^d E[|z_i|^3]$$

for all $t \in \mathbf{R}$ and some constant C .

Clearly we can write $P = \sum_{i=1}^d z_i$ where $z_i = D'_i x_i$ and each D'_i is a random ± 1 . Note that $\sum_{i=1}^d E[|z_i|^3] = \|x\|_3^3$. Let β be the constant $\int_{-\infty}^\infty |t| d\bar{\varphi}(t)$ (the expectation of the absolute value of a Gaussian).

$$\begin{aligned} |E[|P|] - \beta| &= \left| \int_{-\infty}^\infty |t| dF(t) - \int_{-\infty}^\infty |t| d\bar{\varphi}(t) \right| \\ &\leq \int_{-\infty}^\infty |F(t) - \bar{\varphi}(t)| dt \\ &\leq \|x\|_3^3 \int_{-\infty}^\infty \frac{C}{1 + |t|^3} dt. \end{aligned}$$

We claim that $\|x\|_3^3 = O(k^{-1})$. To see this, recall that $\|x\|_2 = 1$, $\|x\|_4 = O(d^{-1/4})$. Equivalently, $\|x^T\|_{2 \rightarrow 2} = 1$

and $\|x^T\|_{4/3 \rightarrow 2} = O(d^{-1/4})$. By applying Riesz-Thorin, we get that $\|x\|_3 = \|x^T\|_{3/2 \rightarrow 2} = O(d^{-1/6})$, hence $\|x\|_3^3 = O(d^{-1/2})$. Since $k = O(d^{1/2})$ the claim is proved.

By linearity of expectation we get $E(W) = k\beta(1 \pm O(k^{-1}))$. We now bound the distance of the median from the expected value.

$$\begin{aligned} |E(W) - \mu| &\leq E[|W - \mu|] \\ &= \int_0^\infty \Pr[|W - \mu| > t] dt \\ &\leq \int_0^\infty 4e^{-t^2/(8\sigma^2)} dt = O(k^{1/2}) \end{aligned}$$

(we used our estimate $\sigma = O(k^{1/2})$ above.) We conclude that $\mu = k\beta(1 + O(k^{-1/2}))$. This clearly shows that (up to normalization) the random transformation $BD\Phi^{(r)}$ (where $r = \lceil 1/\delta \rceil$) has the JL property with respect to embedding into Manhattan space. The running time is $O(d \log d)$.

7 Trimmed Walsh-Hadamard transform

We prove Theorem 2.1. For simplicity, let $H = H_d$. It is well known that computing the Walsh-Hadamard transform $H\mathbf{x}$ requires $O(d \log d)$ operations. It turns out that it is possible to compute $PH\mathbf{x}$ with $O(d \log k)$ operation, as long as the matrix P contains at most k nonzeros. This will imply Theorem 2.1, because code matrices of size $k \times d$ are a product of PH_d , where P contains k rows with exactly one nonzero in each row. To see this we remind the reader that the Walsh-Hadamard matrix (up to normalization) can be recursively described as

$$H_1 = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}, \quad H_q = \begin{pmatrix} H_{q/2} & H_{q/2} \\ H_{q/2} & -H_{q/2} \end{pmatrix}$$

We define \mathbf{x}_1 and \mathbf{x}_2 to be the first and second halves of \mathbf{x} . Similarly, we define P_1 and P_2 as the left and right halves of P respectively.

$$\begin{aligned} (7.9) \quad PH_q \mathbf{x} &= \begin{pmatrix} P_1 & P_2 \end{pmatrix} \begin{pmatrix} H_{q/2} & H_{q/2} \\ H_{q/2} & -H_{q/2} \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \\ &= P_1 H_{q/2} (\mathbf{x}_1 + \mathbf{x}_2) + P_2 H_{q/2} (\mathbf{x}_1 - \mathbf{x}_2) \end{aligned}$$

P_1 and P_2 contain k_1 and k_2 nonzeros respectively, $k_1 + k_2 = k$, giving the recurrence relation $T(d, k) = T(d/2, k_1) + T(d/2, k_2) + d$ for the running time. The base cases are $T(d, 0) = 0$ and $T(d, 1) = d$. We use

induction to show that $T(d, k) \leq 2d \log(k + 1)$.

$$\begin{aligned}
T(d, k) &= T(d/2, k_1) + T(d/2, k_2) + d \\
&\leq d \log(2(k_1 + 1)(k_2 + 1)) \\
&\leq d \log((k_1 + k_2 + 1)^2) \\
&\quad \text{for } k_1 + k_2 = k \geq 1 \\
&\leq 2d \log(k + 1)
\end{aligned}$$

The last sequence of inequalities together with the base cases clearly also give an algorithm and prove Theorem 2.1.

Since in [1] both Hadamard and Fourier transforms were considered we shortly describe also a simple trimmed Fourier transform. In order to compute k coefficients from a d dimensional Fourier transform on a vector \mathbf{x} , we divide \mathbf{x} into L blocks of size d/L and begin with the first step of the Cooley Tukey algorithm which performs d/L FFT's of size L between the blocks (and multiplies them by twiddle factors). In the second step, instead of computing FFT's inside each block, each coefficient is computed directly, by summation, inside it's block. These two steps require $(d/L) \cdot L \log(L)$ and kd/L operations respectively. By choosing $k/\log(k) \leq L \leq k$ we achieve a running time of $O(d \log(k))$.

8 Future work

- *Lower bounds.* A lower bound on the running time of applying a random matrix with a JL property on a vector would be extremely interesting. Any non-trivial (superlinear) bound for the case $k = d^{\Omega(1)}$ will imply a lower bound on the time to compute the Fourier transform, because the bottleneck of our constructions is a Fourier transform.
- *Going beyond $k = d^{1/2-\delta}$.* As part of our work in progress, we are trying to push the result to higher values of the target dimension k (the goal is a running time of $O(d \log d)$). We conjecture that this is possible for $k = d^{1-\delta}$, and have partial results in this direction. A more ambitious goal is $k = \Omega(d)$.

9 Acknowledgements

We thank Tali Kaufman, Bernard Chazelle and Mark W. Tygert for helpful discussions.

References

- [1] N. Ailon and B. Chazelle. Approximate nearest neighbors and the fast Johnson-Lindenstrauss transform. In *Proceedings of the 38st Annual Symposium on the Theory of Computing (STOC)*, pages 557–563, Seattle, WA, 2006.
- [2] N. Alon. Problems and results in extremal combinatorics–I. *Discrete Mathematics*, 273(1-3):31–53, 2003.
- [3] R. I. Arriaga and S. Vempala. An algorithmic theory of learning: Robust concepts and random projection. In *FOCS '99: Proceedings of the 40th Annual Symposium on Foundations of Computer Science*, page 616, Washington, DC, USA, 1999. IEEE Computer Society.
- [4] S. Artstein-Avidan and V. Milman. Logarithmic reduction of the level of randomness in some probabilistic geometric constructions. *SIAM Journal on Computing*, 1(34):67–88, 2004.
- [5] J. Bergh and J. Lofstrom. *Interpolation Spaces*. Springer-Verlag, 1976.
- [6] P. Drineas and R. Kannan. Fast monte-carlo algorithms for approximate matrix multiplication. In *IEEE Symposium on Foundations of Computer Science*, pages 452–459, 2001.
- [7] P. Drineas, R. Kannan, and M. Mahoney. Fast monte carlo algorithms for matrices ii: Computing a low-rank approximation to a matrix, 2004.
- [8] P. Drineas, R. Kannan, and M. Mahoney. Fast monte carlo algorithms for matrices iii: Computing a compressed approximate matrix decomposition, 2004.
- [9] P. Frankl and H. Maehara. The Johnson-Lindenstrauss lemma and the sphericity of some graphs. *Journal of Combinatorial Theory Series A*, 44:355–362, 1987.
- [10] A. M. Frieze, R. Kannan, and S. Vempala. Fast monte-carlo algorithms for finding low-rank approximations. In *IEEE Symposium on Foundations of Computer Science*, pages 370–378, 1998.
- [11] S. Har-Peled. A replacement for Voronoi diagrams of near linear size. In *Proceedings of the 42nd Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 94–103, Las Vegas, Nevada, USA, 2001.
- [12] C. S. Hermann König and N. T. Jaegermann. Projection constants of symmetric spaces and variants of khintchine's inequality. *J. Reine Angew. Math*, 511:1–42, 1999.
- [13] P. Indyk. On approximate nearest neighbors in non-Euclidean spaces. In *Proceedings of the 39th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, pages 148–155, 1998.
- [14] P. Indyk. Uncertainty principles, extractors, and explicit embeddings of ℓ_2 into ℓ_1 . In *Proceedings of the 39th Annual ACM Symposium on the Theory of Computing*, 2007.
- [15] P. Indyk and R. Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In *Proceedings of the 30th Annual ACM Symposium on Theory of Computing (STOC)*, pages 604–613, 1998.
- [16] W. B. Johnson and J. Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. *Contemporary Mathematics*, 26:189–206, 1984.
- [17] S. Khot. Hardness of approximating the shortest

vector problem in lattices. In *Proceedings of the 45th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, 2004.

- [18] E. Kushilevitz, R. Ostrovsky, and Y. Rabani. Efficient search for approximate nearest neighbor in high dimensional spaces. *SIAM Journal on Computing*, 30(2):457–474, 2000.
- [19] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*. Springer-Verlag, 1991.
- [20] N. Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Machine Learning*, 2(4):285–318, 1988.
- [21] F. MacWilliams and N. Sloane. *The Theory of Error Correcting Codes*. North-Holland, 1983.
- [22] J. Matousek. On variants of the Johnson-Lindenstrauss lemma. *Private communication*, 2006.
- [23] A. A. Razborov. Expander codes and somewhat Euclidean sections in ℓ_1^n . *ECCC*, 2007.
- [24] T. Sarlós. Improved approximation algorithms for large matrices via random projections. In *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, Berkeley, CA, 2006.
- [25] M. Sudan. Essential coding theory (class notes).
- [26] S. Vempala. *The Random Projection Method*. DIMACS Series in Discrete Mathematics and Theoretical Computer Science. 2004.
- [27] F. Woolfe, E. Liberty, V. Rokhlin, and M. Tygert. A fast randomized algorithm for the approximation of matrices. *Yale Computer Science Technical Reports, YALE/DCS/TR1380*, 2007.

A Reducing the running time to $O(d \log k)$ for $k \leq d^{1/2-\delta}$

Recall the construction in Section 5: $\delta > 0$ is an arbitrarily small constant, we assume that $k \leq d^{1/2-\delta}$, that k^δ is an integer and that $\beta = f_{\text{BCH}}(k)k^\delta$ divides d (all these requirements can be easily satisfied by slightly reducing δ and at most doubling d). The matrix B is of size $k \times d$, and was defined as follows:

$$B = (B_k \quad B_k \cdots B_k),$$

where B_k is the $k \times f_{\text{BCH}}(k)$ code matrix from Lemma 4.1. Let \hat{B} denote k^δ copies of B_k , side by side. So \hat{B} is of size $k \times \beta$ and B consists of d/β copies of \hat{B} . As in Section 5 we start our construction by studying the distribution of the ℓ_2 estimator $Y = \|BDx\|_2$, where D is our usual random ± 1 diagonal matrix. Going back to (5.2) (recall that M is the matrix whose i 'th column

$M^{(i)}$ is $x_i B^{(i)}$), we recompute the deviation σ :

$$\begin{aligned} \sigma &= \|M\|_{2 \rightarrow 2} = \sup_{\substack{y \in \ell_2^k \\ \|y\|=1}} \|y^T M\|_2 \\ &= \sup \left(\sum_{i=1}^d x_i^2 (y^T B^{(i)})^2 \right)^{1/2} \\ &= \sup \left(\sum_{j=1}^{d/\beta} \sum_{i \in I_j} x_i^2 (y^T B^{(i)})^2 \right)^{1/2}, \end{aligned}$$

where I_j is the j 'th block of β consecutive integers between 1 and d . Applying Cauchy-Schwartz, we get

$$\begin{aligned} (1.10) \quad \sigma &\leq \sup_{\substack{y \in \ell_2^k \\ \|y\|=1}} \left(\sum_{j=1}^{d/\beta} \|x_{I_j}\|_4^2 \|y^T \hat{B}\|_4^2 \right)^{1/2} \\ &= \left(\sup \|y^T \hat{B}\|_4 \right) \|x\|_{(4,2)} = \|\hat{B}^T\|_{2 \rightarrow 4} \|x\|_{(4,2)}, \end{aligned}$$

where $\|\cdot\|_{(p_1, p_2)}$ is defined by

$$\|x\|_{(p_1, p_2)} = \left(\sum_{j=1}^{d/\beta} \|x_{I_j}\|_{p_1}^{p_2} \right)^{1/p_2}$$

and $x_{I_j} \in \ell_{p_1}^\beta$ is the projection of x onto the set of coordinates I_j . Our goal, as in Section 5, is to get $\sigma = O(k^{-1/2})$. By the properties of dual BCH code matrices (Lemma 5.1), we readily have that $\|\hat{B}^T\|_{2 \rightarrow 4} = O((f_{\text{BCH}}(k)k^\delta)^{1/4} k^{-1/2})$ which is $O(k^{\delta/4})$ by our construction. We now need to somehow "ensure" that $\|x\|_{(4,2)} = O(k^{-1/2-\delta/4})$ in order to complete the construction.

As before, we cannot directly control x (and its norms), but we can multiply it by random orthogonal matrices without losing ℓ_2 information. Let H' be a block diagonal $d \times d$ matrix with d/β blocks of the Walsh-Hadamard matrix H_β :

$$H' = \begin{pmatrix} H_\beta & & & \\ & H_\beta & & \\ & & \ddots & \\ & & & H_\beta \end{pmatrix}.$$

Let D' be a random diagonal $d \times d$ matrix over ± 1 . The random matrix $H'D'$ is orthogonal. We study the random variable $X' = \|H'D'x\|_{(4,2)}$. Let M' be the matrix with the i 'th column $M'^{(i)}$ defined as $x_i H'^{(i)}$. We notice that X' is the norm of the Rademacher random variable in $\ell_{(4,2)}^d$ corresponding to M .

Remark: The results on Rademacher random variables, presented in Section 3, apply also to "nonstandard" norms such as $\|\cdot\|_{(p_1, p_2)}$. The dual of $\|\cdot\|_{(p_1, p_2)}$ is $\|\cdot\|_{(q_1, q_2)}$, where q_1, q_2 are the usual dual norm indices of p_1, p_2 , respectively. It is an exercise to check that $\|x\|_{(p_1, p_2)} = \sup_{\|y\|_{(q_1, q_2)}=1} x^T y$. We compute the deviation σ' and a median μ' of X' (as we did in (5.4)):

$$\begin{aligned} \sigma' &= \|M\|_{2 \rightarrow (4,2)} = \|M^T\|_{(4/3,2) \rightarrow 2} \\ &= \sup_{\substack{y \in \ell_{(4/3,2)}^k \\ \|y\|=1}} \left(\sum_i x_i^2 (y^T H^{(i)})^2 \right)^{1/2} \\ &= \sup \left(\sum_{j=1}^{d/\beta} \sum_{i \in I_j} x_i^2 (y^T H^{(i)})^2 \right)^{1/2} \\ &\leq \sup \left(\sum_{j=1}^{d/\beta} \|x_{I_j}\|_4^2 \|y_{I_j}^T H_\beta\|_4^2 \right)^{1/2} \\ &\leq \sup \left(\sum_{j=1}^{d/\beta} \|x_{I_j}\|_4^2 \|y_{I_j}\|_{4/3}^2 \|H_\beta^T\|_{4/3 \rightarrow 4}^2 \right)^{1/2} \\ &= \|H_\beta\|_{4/3 \rightarrow 4} \sup \left(\sum_{j=1}^{d/\beta} \|x_{I_j}\|_4^2 \|y_{I_j}\|_{4/3}^2 \right)^{1/2}, \end{aligned}$$

where the first inequality is Cauchy-Schwartz. By the inequality $(\sum_j A_j)^{1/2} \leq \sum_j A_j^{1/2}$ holding for all nonnegative A_1, A_2, \dots , we get

$$\sigma' \leq \|H_\beta\|_{4/3 \rightarrow 4} \sup_{\substack{y \in \ell_{(4/3,2)}^k \\ \|y\|=1}} \sum_{j=1}^{d/\beta} \|x_{I_j}\|_4 \|y_{I_j}\|_{4/3} \leq \|H_\beta\|_{4/3 \rightarrow 4} \|x\|_{(4,2)}.$$

(The rightmost inequality is from the fact that $\sum_{j=1}^{d/\beta} \|y_{I_j}\|_{4/3}^2 = 1$ and the definition of $\|x\|_{(4,2)}$.) By Hausdorff-Young, $\|H_\beta\|_{4/3 \rightarrow 4} \leq \beta^{-1/4} = O(k^{-1/2-\delta/4})$, hence $\sigma' = O(k^{-1/2-\delta/4} \|x\|_{(4,2)})$. Any median μ' of X' is $O(k^{-1/2-\delta/4})$ (details omitted). Applying Theorem 3.3, we get that for all $t \geq 0$,

$$\Pr[X' > \mu' + t] \leq 4e^{-t^2/(8\sigma'^2)} \leq \hat{c}_1 \exp\{-\hat{c}_2 t^2 k^{1+\delta/2} / \|x\|_{(4,2)}^2\},$$

for some global $\hat{c}_1, \hat{c}_2 > 0$. Setting $t = \Theta(\|x\|_{(4,2)} k^{-\delta/4})$, we get that

$$\Pr[\|H' D' x\|_{(4,2)} > \mu' + t] = O(e^{-k}).$$

Similarly to the arguments leading to Lemma 5.2, and with possible readjustment of the parameter δ , we get using a union bound

LEMMA A.1. [$\ell_{(4,2)}$ **reduction for $k < d^{1/2-\delta}$**] *Let H', D' be as above, and let $\Phi' = H' D'$. Define $\Phi'^{(r)}$ to be a composition of r i.i.d. matrices, each drawn from the same distribution as Φ' . Then With probability $1 - O(e^{-k})$*

$$\|\Phi'^{(r)} x\|_{(4,2)} = O(k^{-1/2-\delta/4})$$

for $r = \lceil 1/2\delta \rceil$.

Combining the above, the random transformation $A = BD\Phi'^{(r)}$ has the *JL* Euclidean property for $k < d^{1/2-\delta}$, and can be applied to a vector in time $O(d \log k)$, as required. Indeed, multiplying by Φ' is done by doing a Walsh transform on d/β blocks of size β each, resulting in time $O(d \log k)$. Clearly the number of random bits used in choosing A is $O(d)$.