

Entropy Matters: Understanding Performance of Sparse Random Embeddings

Maciej Skorski

University of Luxembourg, Luxembourg

Abstract

This work shows how the performance of sparse random embeddings depends on the Renyi entropy-like property of data, improving upon recent works from NIPS'18 and NIPS'19.

While the prior works relied on involved combinatorics, the novel approach is simpler and modular. As the building blocks, it develops the following probabilistic facts of general interest:

- (a) a comparison inequality between the linear and quadratic chaos
- (b) a comparison inequality between heterogenic and homogenic linear chaos
- (c) a simpler proof of Latala's strong result on estimating distributions of IID sums
- (d) sharp bounds for binomial moments in all parameter regimes.

2012 ACM Subject Classification Mathematics of computing → Probabilistic algorithms; Theory of computation → Random projections and metric embeddings

Keywords and phrases Random Embeddings, Sparse Projections, Renyi Entropy

Digital Object Identifier 10.4230/LIPIcs.ISAAC.2022.18

1 Introduction

The celebrated result due to Johnson and Lindenstrauss [38] states that random linear mappings are perfect embeddings: they *nearly preserve distances* of input data points, while mapping them into a *much lower dimension*. This enables accomplishing otherwise computationally demanding tasks, by running on the *reduced yet representative data*. Formally, the lemma states that for any distortion $\epsilon > 0$ and confidence parameter $0 < \delta < 1$, with the embedding dimension $m = \Theta(\log(1/\delta)\epsilon^{-2})$ and the $m \times n$ matrix A sampled from the appropriately scaled normal or Rademacher distribution, for every vector $x \in \mathbb{R}^n$

$$(1 - \epsilon)\|x\|_2 \leq \|Ax\|_2 \leq (1 + \epsilon)\|x\|_2 \quad \text{with probability } 1 - \delta. \quad (1)$$

For modest but practically meaningful distortion and confidence parameters ϵ, δ and large data dimensions n we obtain $m \ll n$, that is a *significant dimension reduction*; on the other hand nearly-preserving distances (up to a relative factor of ϵ) translates into nearly-preserving scalar products and thus the internal data geometry, making it *representative for many tasks*. Indeed, over the years variants of the *Johnson-Lindenstrauss Lemma* have found important applications to text mining and image processing [7], approximate nearest neighbor search [35, 3], learning mixtures of Gaussians [22], sketching and streaming algorithms [43, 47], approximation algorithms for clustering high dimensional data [6, 12, 54], speeding up linear algebraic computations [57, 61, 16], analyzing combinatorial properties of graphs [28, 52] and even to privacy [9, 42]; on the pure theory side, it is worth mentioning the importance for understanding Hilbert spaces in functional analysis [39].

Although the embedding dimension m is optimal [40, 37], the costly matrix-vector product can be optimized by the use of *sparse matrices*. The long line of research [1, 21, 51, 3, 55, 41, 18] have finally established the same guarantees for matrices A with only $s = \Theta(\log(1/\delta)\epsilon^{-1})$



© Maciej Skorski;

licensed under Creative Commons License CC-BY 4.0

33rd International Symposium on Algorithms and Computation (ISAAC 2022).

Editors: Sang Won Bae and Heejin Park; Article No. 18; pp. 18:1–18:18

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

entries per column¹. Optimal in the worst-case, these results were far away from the performance empirically observed on real-world data, particularly the remarkable accuracy of *feature hashing* [65] which uses only $s = 1$ (!). This led to the following intriguing question:

Why extremely sparse random projections work better-than-expected?

A careful reader notices that so far we have been speaking of *data-oblivious* results, that is under no assumption on the data structure. Indeed, the relevant research in [65, 21, 41, 30, 36] has finally established [30, 36] that the certain metric which captures *data dispersion*, more precisely the ratio $v = \|x\|_\infty / \|x\|_2$, allows for setting the matrix sparsity to²

$$s = \Theta(v^2 \epsilon^{-1}) \cdot \max \left\{ \log \frac{1}{\delta}, \frac{\log^2 \frac{1}{\delta}}{\log^2 \frac{1}{\epsilon}} \right\} \quad (2)$$

while keeping the optimal dimension $m = \Theta(\log(1/\delta)\epsilon^{-2})$. This offers an additional improvement by a factor of $1/v^2$. In simple terms: the more data is dispersed, the better matrix sparsity works. This breakthrough result still suffers from the following limitations:

1. *Unsatisfactory definition of data dispersion.* The ratio ℓ_∞ -to- ℓ_2 is a crude notion: on the unit sphere $\|x\|_2 = 1$ it depends on the heaviest element and so is not smooth enough. It suffers particularly from “spikes” that are naturally present in real-world data (such as features produced in text-mining [4]) and due to pairwise vector differences studied in multi-vector setup (uniform guarantees for multiple vectors are obtained by looking at pairwise differences $x - x'$, which leads to “spikes” for example in images [48]). This motivates further research for a *more accurate notion of dispersion*.
2. *High proof complexity and lack of modern toolkit.* Proofs in prior works [30, 36] suffer from being lengthy and convoluted, mostly in supplementary materials, which results in numerical mistakes as well as gaps not immediately fixable (see Appendix A). These works did admirable efforts on presenting the self-contained proof, yet did not utilize the modern probability toolkit to the full extent. Their strategy is to see Equation (1) as the concentration of the *quadratic form* $x \rightarrow \|Ax\|_2^2$, and quantify its tails by controlling high-order moments estimated via multinomial expansions coupled with combinatorial arguments. However, this does not leverage tools to control quadratic random forms, namely the modern techniques of the *Hanson-Wright inequality* [31, 60, 67] such as decoupling of quadratic forms [63, 24]. Furthermore, it re-develops a variant of the sharp result from [49] on moment estimation and certain known facts from *high-dimensional probability* on sub-gaussian distributions [11, 10]. Finally, while [36] develops its technical lemmas for symmetric random variables, this condition is not satisfied which leaves a gap. Thus, further effort in *revisiting and modernizing the toolkit* used in recent state-of-art works [30, 36] is well-motivated. Indeed, simplifying proofs and developing novel techniques for the JL Lemma is an independent and valued line of research [28, 29, 23, 19], as these have been historically difficult (the original result used sophisticated geometric approximations, while the sparse variant [21] relied on correlation inequalities [27]).

¹ As shown by [18] one can reduce further sparsity s by $B > 1$ at the cost of exponentially increasing the dimension m by a factor of $2^{\Theta(B)}$. However, in practice, sub-optimal dimension is less interesting.

² The formula arises from rearranging Theorem 1.5 in [36]

2 Our Contribution

This work offers a solution to the two problems discussed above: we *strengthen* and to a great extent *simplify* the state-of-art results from prior works.

2.1 Performance of Sparse Random Projections

We introduce the following (novel) notion of the *data dispersion*:

$$v_d(x) \triangleq \sup_{|I| < d/2} \left(\frac{\sum_{i \notin I} |x|_i^d}{\sum_{i \notin I} x_i^2} \right)^{\frac{1}{d-2}} / \|x\|_2, \quad d > 2. \quad (3)$$

where I are taken as strict subsets of the support of x .

The matrix A is sampled from the sparsified Rademacher distribution, as in prior works:

Algorithm 1 Sparse Random Projections: Matrix Sampler.

Data: data dimension n , embedding dimension m , matrix sparsity s

Result: $A \in \mathbb{R}^{n \times m}$

- for every column i , select s positions at random (without replacement)
- set randomly ± 1 on the selected positions
- scale the matrix by $1/\sqrt{s}$

For the matrix as in Algorithm 1 above, we prove the following result.

► **Theorem 1.** *Let $d = \log(1/\delta)$, then the JL Lemma, that is (1), holds for the dimension*

$$m = \Theta(d\epsilon^{-2}) \quad (4)$$

and any sparsity s such that

$$v_d(x) \leq \Theta(s\epsilon)^{1/2} \min(\log(m\epsilon/d)/d, 1/d^{1/2}). \quad (5)$$

We now discuss the result in detail in the series of remarks below.

► **Remark 2 (Intuition).** We give the following rationale for one could conjecture a result like the one above: the analysis of sparse random projections establishes that the performance depends on the d -th moment of the error expression, where $d = \log(1/\delta)$ is relatively small; it seems reasonable to expect that the assumptions on the data should not include moments higher than of order d , particularly bounding $\|x\|_\infty$ seems to be an overshooting.

► **Remark 3 (Comparison with previous bounds).** Since $v_d(x) \leq \|x\|_\infty / \|x\|_2$, we obtain the previous state-of-art bounds from [36], by rearranging Equation (5) to Equation (2). This approximation is however rather crude, as it merely replaces the d -th norm $\|\cdot\|_d$ by $\|\cdot\|_\infty$, and our bound can do much better. Consider the more explicit example where $x_i^2 = (n/d)^{-1/d}$ for d values of i and $x_i^2 = 1 - (n/d)^{-1/d}/(n-d)$ otherwise. We then have $v_d(x) = \Theta(n^{-\frac{2}{d-2}})$ while $\|x\|_\infty / \|x\|_2 = \Theta(n^{-\frac{1}{d}})$. Since the best possible sparsity s is roughly proportional to $v_d(x)^{-2}$, our gain over the previous approach is by a factor of $n^{\frac{4}{d-2} - \frac{2}{d}}$ which is huge for moderate values of d and large n (that is, in a typical application regime).

► **Remark 4 (Relation to Renyi Entropy).** Let's introduce the probability measure $w_i \sim x_i^2$, then $(\sum_i |x_i|^d / \sum_i x_i^2)^{\frac{1}{d-2}} / \|x\|_2 = (\sum_i w_i^{\frac{d}{2}})^{\frac{1}{d-2}} = 2^{H_{d/2}((w_i))/2}$ where the Renyi entropy [58] of the distribution w is defined as $H_d(w) \triangleq \frac{1}{1-d} \sum_i w_i^d$ and $H_\infty(w) \triangleq -\log \max_i w_i$ when

$d = \infty$. Under the mild assumption that x such that $\sum_{i \notin I} x_i^2 = \Theta(\|x\|_2^2)$ for all $|I| \leq d$ we can thus compare the sparsity achieved in Theorem 1 and the result in [36] as low-order Renyi entropy versus min-entropy. More precisely, our bound on s is better by a factor of $2^{H_{d/2}((w_i)) - H_\infty((w_i))}$, that is the gain is *exponential in entropy deficiency* understood as $H_{d/2}((w_i)) - H_\infty((w_i))$. The well-known bounds from information-theory [14] show that this gap can be as big as $\frac{1}{d/2-1} H_{d/2}((w_i))$ (which is unbounded without some restrictions on x).

► **Remark 5 (Dimension-Sparsity Tradeoffs).** It is possible to improve the sparsity parameter s by a factor of B at the expense of making the dimension worse by a factor of $e^{\Theta(B)}$, exactly as in [36]. However, this tradeoff does not seem to be interesting from the application-oriented point of view (the whole idea of random projections is to keep the low dimension).

2.2 Techniques of Independent Interest

2.2.1 From Quadratic to Linear Chaos

One important novelty in our approach is that we get rid of analyzing quadratic forms, which appear due to considering the expression $\|Ax\|_2^2$, by an elegant reduction to their linear analogues. Although quadratic chaoses of symmetric random variables have been studied in the past [49, 46], the generic bounds were found intractable to analyze by the authors of prior works [30, 36] and other workarounds have been proposed. It has been not clear if one can get rid of these complicated methods. Indeed, we show that we can:

► **Lemma 6.** *Let X_i be independent zero-mean random variables, with possibly different distributions. Then for even $d \geq 2$ we have*

$$\left\| \sum_{i \neq j} X_i X_j \right\|_d \leq 32 \left\| \sum_i X_i \right\|_d^2.$$

► **Remark 7.** The result is fairly general, not requiring symmetry or identical distributions. In fact, the constant reduces to 4 if X_i are already symmetric.

This bound allows for reducing a bulk of technical calculations, and almost directly applying existing *tractable bounds* for linear forms such as those in [50]. The proof uses *decoupling* [63] which allows for upper-bounding the moments of the quadratic form $\sum_{i \neq j} X_i X_j$ by the moments of bilinear form $\sum_{i \neq j} X_i X'_j$, and *symmetrization* [64] which allows for replacing X_i by their symmetrized versions $X_i - X'_i$ at the expense of a constant factor.

2.2.1.1 Heterogenic Sparse Rademacher Chaos

Although we reduce the problem to studying linear forms, they are not IID sums. More precisely in our case we will be interested in sums of form $\sum_i x_i X_i$ where X_i are symmetric and IID, but the given weights x_i can be very different. Such sums are notoriously difficult to analyze, the best example being probably the classical Khintchine's inequality which seeks to bound $\left\| \sum_i x_i \sigma_i \right\|_d$ where σ_i are Rademachers, for a given sequence of weights (x_i) ; it took a while until the original bounds [44] have been tightened, in a way that explicitly depend on x [33]. While prior works [30, 36] handle this difficulty in our context implicitly (in combinatorial analyses of multinomial expansions), we use *majorization theory* to essentially compare the heterogenic and homogenic (easier) setup. We prove

► **Lemma 8.** Let $\|x\|_2 = 1$ and $X_i \sim^{IID} \eta_i \sigma_i$ where η_i are IID Bernoulli and σ_i are IID Rademacher r.vs. Then for $v = v_d(x)$ where $v_d(x)$ is as in Equation (3), and even $d > 0$

$$\left\| \sum_i x_i X_i \right\|_d \leq O\left(\|K^{-1/2} \sum_{i=1}^K X_i\|_d\right), \quad K = \lceil v^{-2} \rceil.$$

The result depends on the structure of x captured by $v = v_d(x)$, note that the equality holds when $x_i = v$ for all non-zero weights x_i (note that we normalize $\|x\|_2 = 1$ w.l.o.g.); this is the core of our method, and we can see it as a sparse analogue of Khintchine's Inequality (Bernoulli variables restrict the summation to a random subset). The result should be considered strong and somewhat surprising; per analogy to the case when there are no Bernoulli variables, results from majorization theory seem to suggest that the moment should be rather minimized for x_i that are nearly uniform³. The answer is in the condition $v_d(x)$ which is, to a certain degree, a relaxation of the requirement that x_i is flat and in the constant under $O(1)$. What we prove is not that (x_i) with K elements gives the maximum, but that the value differs from the actual maximum by at most a constant factor. In our proof, we use the assumption in Equation (3) and majorization [17] to compare the behavior of sums $S_k = \sum_{i_1 \neq \dots \neq i_k} x_{i_1}^2 \cdots x_{i_k}^2$ when x_i is uniform over K elements versus over the whole space. Under the normalizing condition $\|x\|_2 = 1$, they can be interpreted as *birthday collision probabilities*, which makes the comparison easy to evaluate.

2.2.1.2 Moments of IID Sums

We will need a result which provides *tight bounds on moments of iid sums*. Although this problem has been solved by a characterization due to Latala [50], the result seems to be little known within the TCS community; instead classical bounds due to Hoeffding [34], Chernoff [15], Bernstein [5] or more modern bounds stated sub-gaussian or sub-gamma distributions [11] are used. Since the analysis of sparse random projections involves random variable with little exotic behavior, the classical inequalities are not sufficient.

In hope for popularizing the technique and to make the paper self-consistent, we provide an alternative and simpler proof of Latala's result [50].

► **Lemma 9.** For zero-mean r.vs. $X_i \sim^{IID} X$ and even $d > 0$

$$\left\| \sum_{i=1}^n X_i \right\|_d \leq 2e \cdot \max_k \left[\binom{d}{k}^{1/k} (\exp(d/n) - 1)^{-1/k} \|X\|_k : \max(2, d/n) \leq k \leq d \right] \quad (6)$$

which implies the following simpler bound

$$\left\| \sum_{i=1}^n X_i \right\|_d \leq \frac{2e^2}{(1 - e^{-1})^{1/2}} \cdot \max_k \left[d/k \cdot (n/d)^{1/k} \cdot \|X\|_k : \max(2, d/n) \leq k \leq d \right]. \quad (7)$$

► **Remark 10.** In addition to simplifying the proof, we provide an explicit constant (not given in the original proof). For non-symmetric distributions, our numerical constant is better than the one implied by symmetrizing the original proof. We also note that there is the same matching, up to a constant, lower bound [49], so that in the result above we have the equality up to a constant.

³ The map $(x_i) \rightarrow \left\| \sum_i x_i \sigma_i \right\|_d^d$ is Schur-concave in variables x_i^2 [26].

2.2.1.3 Sharp Bounds for Binomial Moments

Having reduced the problem to studying moments of $\sum_i \eta_i \sigma_i$, we face the problem of estimating binomial moments. Somewhat surprisingly, the literature does not offer good bounds for binomial moments. What we know are combinatorial formulas [45] not in a closed asymptotic form, and nearly perfect estimates (up to $o(1)$ relative error) for binomial probabilities [62] as well as the tails [20, 53, 56] (see also the survey in [2]); these tails unfortunately lead to intractable integrals expressing moments (with Kullback-Leibler terms).

Since the question is foundational with clear potential for applications beyond our problem, we give the following general and detailed answer

► **Lemma 11.** *Let $S \sim \text{Binom}(K, p)$ where $p \leq \frac{1}{2}$, and $d > 0$ be even. Then*

$$\|S - \mathbf{E}S\|_d = \Theta(1) \begin{cases} (dKp)^{1/2} & \log(d/Kp) < d/K \leq 2 \\ Kp^{K/d} & \log(d/Kp) < 2 \leq d/K \\ \frac{d}{\log(d/Kp)} & \max(2, d/K) \leq \log(d/Kp) \leq d \\ (Kp)^{1/d} & d < \log(d/Kp) \end{cases} \quad (8)$$

► **Remark 12.** The bound has up to 4 regimes, in which we provide an estimate sharp up to a constant. The upper bound (sufficient for our needs) follows from Lemma 9, while the lower bound holds because the bound in Lemma 9 is sharp up to an absolute constant [49].

2.3 Proof Outline

We actually prove that

$$(1 - \epsilon)\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \epsilon)\|x\|_2^2 \quad \text{with probability } 1 - \delta \quad (9)$$

from which Equation (1) follows by taking the square roots and using the elementary inequalities $\sqrt{1 + \epsilon} \leq 1 + \epsilon$, $1 - \epsilon \leq \sqrt{1 - \epsilon}$. Denoting $Z = \|Ax\|_2^2 - \|x\|_2^2$ we find [36])

$$Z = \frac{1}{s} \sum_{r=1}^m Z_r, \quad Z_r \triangleq \sum_{i \neq j} x_i x_j \eta_i \eta_j \sigma_i \sigma_j. \quad (10)$$

It can be shown that Z_r are *negatively dependent* and thus their sum obey moment upper-bounds for independent random variables [25, 8]. More precisely we have that

$$\|Z\|_d \leq \frac{1}{s} \left\| \sum_{r=1}^m Z_r \right\|_d, \quad Z_r \sim^{IID} \sum_{i \neq j} x_i x_j \eta_i \eta_j \sigma_i \sigma_j. \quad (11)$$

The techniques outlined above, namely Lemma 6 and Lemma 8 show that for $K = \lceil v_d(x)^{-2} \rceil$

$$\|Z_r\|_d \leq O(K^{-1} \|S - S'\|_d^2), \quad S, S' \sim^{IID} \text{Binom}(K, p). \quad (12)$$

Since $\|S - S'\|_d \leq 2\|S - \mathbf{E}S\|_d$ (the triangle inequality), by Lemma 11 we obtain

► **Corollary 13.** *For any even $d > 0$ we have*

$$\|Z_r\|_d \leq O(1) \begin{cases} dp & \log(d/Kp) < d/K \leq 2 \\ Kp^{2K/d} & \log(d/Kp) < 2 \leq d/K \\ \frac{K^{-1}d^2}{\log^2(d/Kp)} & \max(2, d/K) \leq \log(d/Kp) \leq d \\ K^{-1}(Kp)^{2/d} & d < \log(d/Kp) \end{cases} \quad (13)$$

It now suffices to plug this bound in Lemma 8 (it applies for negatively dependent r.v.s.) and analyze the 4 different regimes, to obtain moment bounds for Z from Equation (10); then Theorem 1 follows by Markov's inequality. The work has been mostly finalized at this point, due to our modular approach; the application of Lemma 8 is discussed in the appendix.

► **Remark 14.** At the final stage [36] also obtains analogous bounds (with K defined in terms of $v = \|x\|_\infty/\|x\|_2$). They are however not derived via a single application of a lemma, but rather a mixture of three techniques (direct bounds on quadratic forms, linear forms, and the reproved result on the sub-gaussian norm of a binary random variable [13]).

2.4 Organization

The rest of the paper is organized as follows: in Section 3 we introduce basic notation and some simple auxiliary facts that will be used throughout the discussion, in Section 4 we present proofs of the key ingredients of our proof. Details omitted in the proof outline are provided in Appendix B. In Section 5 we conclude the work.

3 Preliminaries

3.1 Basic Notation

For a random variable X , we define its d -th moment as $\mathbf{E}|X|^d$ and its d -th norm as $\|X\|_d = (\mathbf{E}|X|^d)^{1/d}$ (this is indeed a norm when $d \geq 1$). For the sequence (x_i) we define $\|(x_i)\|_d = (\sum_i |x_i|^d)^{1/d}$ for $0 < d < 1$, $\|x\|_\infty = \max_i |x_i|$ and $\|x_i\|_0 = \#\{i : x_i \neq 0\}$.

By $\text{Bern}(p)$ we denote the Bernoulli distribution, that is 1 with probability p and zero otherwise. By $\text{Binom}(K, p)$ we denote the binomial distribution with parameters K and p (equal in the distribution to the sum of K independent copies of $\text{Bern}(p)$).

3.2 Auxiliary Functions

We need the elementary properties of the two functions that often appear in our analysis:

► **Proposition 15.** *The function $g(d) = 1/q \cdot a^{1/q}$ for $q > 0$ is decreasing when $a \geq 1$ and for $a < 1$ it achieves its local maximum at $q = \log(1/a)$ with the value $g(q) = 1/e \log(1/a)$.*

► **Proposition 16.** *The function $g(q) = q \cdot a^{1/q}$ for $q > 0$ is increasing when $a \leq 1$ and for $a > 1$ achieves its local minimum at $q = \log a$ with the value $g(q) = e \log a$.*

3.3 Probabilistic Techniques

The following fact will allow us to handle non-symmetric distributions.

► **Proposition 17** (Symmetrization trick [64]). *For any norm $\|\cdot\|$ we have*

$$\frac{1}{2} \left\| \sum_i X_i \sigma_i \right\| \leq \left\| \sum_i X_i \sigma_i \right\| \leq 2 \left\| \sum_i X_i \sigma_i \right\|,$$

for any zero-mean independent X_i and independent Rademacher random variables σ_i .

We will also need the decoupling inequality, useful in attacking quadratic forms

► **Proposition 18** (Decoupling inequality [63]). *Let X_i be zero-mean independent r.v.s. and X'_i be their independent copies. Then for any weights $a_{i,j}$*

$$\mathbf{E}f\left(\sum_{i \neq j} a_{i,j} X_i X_j\right) \leq \mathbf{E}f\left(4 \sum_{i \neq j} a_{i,j} X_i X'_j\right),$$

for any convex function f .

► Remark 19. The summation is over $i \neq j$, e.g. the quadratic form must be off-diagonal!

4 Proofs

4.1 Quadratic vs Linear Chaos

Proof of Lemma 6. Let X'_i be independent copies of X_i . The decoupling inequality gives

$$\left\| \sum_{i \neq j} X_i X_j \right\|_d \leq 4 \left\| \sum_{i \neq j} X_i X'_j \right\|_d. \quad (14)$$

We apply the symmetrization trick to the d -th norm twice: first for random variables X_i with any fixed choice of X'_j which gives $\left\| \sum_{i \neq j} X_i X'_j \right\|_d \leq 2 \left\| \sum_{i \neq j} X_i \sigma_i X'_j \right\|_d$ (here we use the independence of X_i and X'_j) and second for random variables X'_j under the fixed values of $X_i \sigma_i$ which gives $\left\| \sum_{i \neq j} X_i X'_j \right\|_d \leq 4 \left\| \sum_{i \neq j} X_i \sigma_i X'_j \sigma'_j \right\|_d$ (σ'_j is an independent Rademacher sequence). For simplicity, we denote $X_i := X_i \sigma_i$ and $X_j := X_j \sigma'_j$, note that the introduced random variables $X_i \sigma_i$ and $X_j \sigma'_j$ are also identically distributed.

Consider the sum $\sum_{i,j} X_i X'_j = \sum_i (\sum_{j \neq i} X'_j) X_i$ as linear in X_i with coefficients depending on X'_j , and apply the multinomial theorem which gives

$$\mathbf{E}\left[\left(\sum_{i \neq j} X_i X'_j\right)^d \middle| (X'_j)\right] = \sum_{(d_i)} \binom{d}{2d_1 \dots 2d_n} \prod_i \left(\sum_{j \neq i} X'_j\right)^{2d_i} \mathbf{E} X_i^{2d_i}.$$

where we use the symmetry of X_i , so that all odd moments vanish. Again by the multinomial theorem we see that

$$\mathbf{E}\left(\sum_{j \neq i} X'_j\right)^d \leq \mathbf{E}\left(\sum_j X'_j\right)^d.$$

Combining the last two bounds gives

$$\begin{aligned} \mathbf{E}\left(\sum_{i \neq j} X_i X'_j\right)^d &\leq \mathbf{E}_{(X'_j)} \left[\mathbf{E}\left[\left(\sum_{i \neq j} X_i X'_j\right)^d \middle| (X'_j)\right] \right] \\ &\leq \sum_{(d_i)} \binom{d}{2d_1 \dots 2d_n} \mathbf{E}\left[\prod_i \left(\sum_{j \neq i} X'_j\right)^{2d_i} X_i^{2d_i}\right] \\ &\leq \mathbf{E}\left(\sum_i \left(\sum_j X'_j\right) X_i\right)^d \\ &= \mathbf{E}\left(\sum_i X_i\right)^d \left(\sum_j X'_j\right)^d = \mathbf{E}\left(\sum_i X_i\right)^{2d}, \end{aligned}$$

which can be stated as

$$\left\| \sum_{i \neq j} X_i X'_j \right\|_d \leq \left\| \sum_i X_i \right\|_d^2. \quad (15)$$

By combining Equation (14) and Equation (15), and keeping in mind that X_i above are symmetrized versions, we obtain for original (only centered) random variables X_i

$$\mathbf{E} \left\| \sum_{j \neq i} X_i X_j \right\|_d \leq 16 \mathbf{E} \left\| \sum_{j \neq i} X_i \sigma_j \right\|_d,$$

and the result follows by one more application of the symmetrization trick. \blacktriangleleft

4.2 Heterogenic vs Homogenic Chaos

Proof of Lemma 8. By the multinomial expansion and the symmetry of Z_i (which implies that the odd moments vanish) we obtain

$$\mathbf{E} \left(\sum_i x_i X_i \right)^d = \sum_{(d_i)} \binom{d}{2d_1 \dots 2d_n} p^{\|(d_i)\|_0} \prod_i x_i^{2d_i},$$

where the summation is over non-negative sequences (d_i) for $i = 1, \dots, n$ such that $\sum_i d_i = d/2$, and we denote $\|(d_i)\|_0 = \#\{i : d_i > 0\}$. Considering possible values of $k = \|(d_i)\|_0$, we find that the above expression is a non-negative combination of

$$S_k^{[d]}(x) = \sum_{i_1 \neq \dots \neq i_k} x_{i_1}^{2d_1} \dots x_{i_k}^{2d_k}$$

where possible values of k are $1 \leq k \leq \min(d/2, n_0)$ where $n_0 = \|(x_i)\|_0$. We now apply our assumption on x iteratively to $x_{i_k}, x_{i_{k-1}}, \dots$, obtaining

$$S_k^{[d]}(x) \leq v^{2 \sum_{i: d_i > 1} (d_i - 1)} \sum_{i_1 \neq \dots \neq i_k} x_{i_1}^2 \dots x_{i_k}^2.$$

Here we have used the fact that $v_d(x)$ is increasing in d , so $v_k(x) \leq v$ when $k \leq d$; this follows from seeing $v_d(x)$ as the power mean of order $d - 2$ and weights $x_i^2 / \sum_{i \notin I} x_i^2$ [32, 66].

We make the following important observation: the equality holds whenever x_i is flat with the value v , e.g. all non-zero entries are equal to v . Observe that the sums $S_k(x) = \sum_{i_1 \neq \dots \neq i_k} x_{i_1}^2 \dots x_{i_k}^2$ are elementary symmetric polynomials in variables $y_i = x_i^2$ where $\sum_i y_i = \sum_i x_i^2 = 1$, hence over the probability simplex. The elementary symmetric functions are Schur-concave [17], and thus they are maximized at the uniform distribution, in our case when $x_i = n^{-1/2}$. In fact, $S_k(x)$ is the probability that k independent samples from the distribution $p_i = x_i^2$ do not collide. For any sequence (x_i^2) which has N non-zero equal entries and $\sum_i x_i^2 = 1$ we have that:

$$S_k(x) = N \cdot (N - 1) \cdot \dots \cdot (N - k + 1) / N^k.$$

Since $N \geq k$ and since $k \leq d$, using Stirling's approximation [59] we obtain

$$S_k(x) = \prod_{i=0}^{k-1} (1 - i/N) \geq k! / k^k = \Theta(1)^k \geq \Theta(1)^d.$$

Clearly $S_k(x) \leq 1$ for any x . If we replace (x_i) by a sequence such that $x_i = v$ for $K = v^{-2}$ values of i (e.g., flat), we lose at most a factor of $\Theta(1)^k \leq \Theta(1)^d$ in every term $S_k^{[d]}(x)$. \blacktriangleleft

4.3 Moments of IID Sums

Proof of Lemma 9. We have the following chain of estimates

$$\begin{aligned}
\mathbf{E}\left(\sum_i X_i\right)^d &= \sum_{d_i: d_1+\dots+d_n=d, d_i \geq 2} \binom{d}{d_1 \dots d_n} \prod_i \mathbf{E} X_i^{d_i} \\
&\leq \sum_{d_i: d_1+\dots+d_n=d, d_i \geq 2} \prod_i \binom{d}{d_i} \mathbf{E} |X_i|^{d_i} \\
&\leq \sum_{d_i \geq 2} \prod_i \binom{d}{d_i} \mathbf{E} |X_i|^{d_i} \\
&\leq \left(\sum_{k=2}^d \binom{d}{k} \|X\|_k^k \right)^n.
\end{aligned}$$

Applying this for $X_i := X_i/t$ we have for any $t > 0$

$$\mathbf{E}\left(t^{-1} \sum_i X_i\right)^d \leq \left(\sum_{k=2}^d \binom{d}{k} \|X\|_k^k / t^k \right)^n.$$

Thus $\|\sum_i X_i\|_d \leq et$ for any t such that the right-hand side is at most e , equivalently

$$\sum_{k=2}^d \binom{d}{k} \|X\|_k^k / t^k \leq \exp(d/n) - 1,$$

which is satisfied for

$$t = 2 \max_{k=2 \dots d} \binom{d}{k}^{1/k} (\exp(d/n) - 1)^{-1/k} \|X\|_k.$$

This proves the first part. Observe that for $k \geq 2$ we have

$$\binom{d}{k}^{1/k} (\exp(d/n) - 1)^{-1/k} \leq \frac{ed}{k \exp(d/kn)} \cdot \frac{1}{(1 - \exp(-1))^{1/2}},$$

where we use the elementary inequalities $\binom{d}{k} \leq (de/k)^k$ and $\exp(u) - 1 \geq \exp(u) \cdot (1 - e^{-1})$ for $u \geq 1$. The function $u \rightarrow u/\exp(u)$ decreases for $u \geq 1$; applying this to $u = d/kn$ gives

$$\binom{d}{k}^{1/k} (\exp(d/n) - 1)^{-1/k} \leq \frac{en}{(1 - e^{-1})^{1/2}}, \quad k \leq d/n.$$

Since $\|X\|_k$ increases in k we have

$$\max_{k=2 \dots d, k \leq d/n} \binom{d}{k}^{1/k} (\exp(d/n) - 1)^{-1/k} \|X\|_k \leq \frac{en \|X\|_{d/n}}{(1 - e^{-1})^{1/2}}.$$

We have $(\exp(d/n) - 1)^{-1/k} \leq (d/n)^{-1/k}$ due to the elementary inequality $\exp(u) - 1 \geq u$, and $\binom{d}{k} \leq (de/k)^k$ for any k . This gives

$$\max_{k=2 \dots d} \binom{d}{k}^{1/k} (\exp(d/n) - 1)^{-1/k} \|X\|_k \leq e \max_{k=2 \dots d} d/k \cdot (n/d)^{1/k} \cdot \|X\|_k$$

When $d/n \geq 2$ we have that $d/k \cdot (n/d)^{1/k} \cdot \|X\|_k = n \|X\|_{d/n} \cdot 2^{-1/2}$ for $k = d/n$. Comparing the last two equations, we obtain

$$\max_{k=2 \dots d, k \leq d/n} \binom{d}{k}^{1/k} (\exp(d/n) - 1)^{-1/k} \|X\|_k \leq C \max_{k=2 \dots d, k > d/n} d/k \cdot (n/d)^{1/k} \cdot \|X\|_k,$$

with $C = \frac{e}{(1 - e^{-1})^{1/2}}$. This completes the proof. \blacktriangleleft

4.4 Binomial Moments

Proof of Lemma 11. Applying Lemma 9 we obtain

$$\|S - \mathbf{E}S\|_d \leq O(1) \cdot \max \left\{ (d/k) \cdot (Kp/d)^{1/k} : \max(2, d/K) \leq k \leq d \right\}.$$

because $S \sim \sum_i X_i$ where $X_i \sim \text{Bern}(p)$ and $\|X_i - \mathbf{E}X_i\|_d = (p(1-p)^{d-1} + (1-p)p^{d-1})^{1/d}$ so that $\|X_i - \mathbf{E}X_i\|_d = \Theta(p)^{1/d}$ for $p \leq 1/2$.

The expression under the maximum is proportional to $k^{-1} \cdot a^{1/k}$ where $a = Kp/d$. The claim follows by applying Proposition 15, namely a) when $\max(2, d/K) \leq \log(1/a) \leq d$ (that is, inside the interval) we have necessarily $a \leq e^{-2} < 1$ our maximum is at $k = \log(1/a)$ b) when $\log(1/a) > d$ we must have $a < 1$ and our maximum is at $k = d$ and c) when $\log(1/a) < \max(2, d/K)$ then the maximum is at $k = \max(2, d/K)$. ◀

5 Conclusion

We have proven novel bounds for sparse random projections, showing that the performance depends on the data statistic closed to *Renyi entropy*. Some intriguing problems we leave for future work are

- How do results extend to non-Rademacher matrices?
- Can we use majorization theory to fully characterize worst case for the linear chaos?

References

- 1 Dimitris Achlioptas. Database-friendly random projections. In *Proceedings of the twentieth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 274–281, 2001.
- 2 Thomas D Ahle. Asymptotic tail bound and applications, 2017.
- 3 Nir Ailon and Bernard Chazelle. Approximate nearest neighbors and the fast Johnson–Lindenstrauss transform. In *Proceedings of the thirty-eighth annual ACM symposium on Theory of computing*, pages 557–563, 2006.
- 4 Akiko Aizawa. An information-theoretic perspective of tf-idf measures. *Information Processing & Management*, 39(1):45–65, 2003.
- 5 SN Bernshtein. Probability theory (in Russian). *Gosizdat, Moscow-Leningrad*, 1927.
- 6 Gérard Biau, Luc Devroye, and Gábor Lugosi. On the performance of clustering in Hilbert spaces. *IEEE Transactions on Information Theory*, 54(2):781–790, 2008.
- 7 Ella Bingham and Heikki Mannila. Random projection in dimensionality reduction: applications to image and text data. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 245–250, 2001.
- 8 Henry W Block, Thomas H Savits, Moshe Shaked, et al. Some concepts of negative dependence. *The Annals of Probability*, 10(3):765–772, 1982.
- 9 Jeremiah Blocki, Avrim Blum, Anupam Datta, and Or Sheffet. The Johnson–Lindenstrauss transform itself preserves differential privacy. In *2012 IEEE 53rd Annual Symposium on Foundations of Computer Science*, pages 410–419. IEEE, 2012.
- 10 Stéphane Boucheron, Olivier Bousquet, Gábor Lugosi, Pascal Massart, et al. Moment inequalities for functions of independent random variables. *The Annals of Probability*, 33(2):514–560, 2005.
- 11 Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- 12 Christos Boutsidis, Anastasios Zouzias, and Petros Drineas. Random projections for k -means clustering. In *Advances in Neural Information Processing Systems*, pages 298–306, 2010.

- 13 V Buldygin and K Moskvichova. The sub-gaussian norm of a binary random variable. *Theory of probability and mathematical statistics*, 86:33–49, 2013.
- 14 Christian Cachin. Smooth entropy and rényi entropy. In *International Conference on the Theory and Applications of Cryptographic Techniques*, pages 193–208. Springer, 1997. URL: https://link.springer.com/chapter/10.1007/3-540-69053-0_14.
- 15 Herman Chernoff et al. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *The Annals of Mathematical Statistics*, 23(4):493–507, 1952.
- 16 Kenneth L Clarkson and David P Woodruff. Low-rank approximation and regression in input sparsity time. *Journal of the ACM (JACM)*, 63(6):1–45, 2017.
- 17 M Lawrence Clevenston and William Watkins. Majorization and the birthday inequality. *Mathematics Magazine*, 64(3):183–188, 1991.
- 18 Michael B Cohen. Nearly tight oblivious subspace embeddings by trace inequalities. In *Proceedings of the twenty-seventh annual ACM-SIAM symposium on Discrete algorithms*, pages 278–287. SIAM, 2016.
- 19 Michael B Cohen, TS Jayram, and Jelani Nelson. Simple analyses of the sparse Johnson–Lindenstrauss transform. In *1st Symposium on Simplicity in Algorithms (SOSA 2018)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2018.
- 20 Harald Cramér. On a new limit theorem of the theory of probability. *Uspekhi Matematicheskikh Nauk*, 10:166–178, 1944.
- 21 Anirban Dasgupta, Ravi Kumar, and Tamás Sarlós. A sparse Johnson-Lindenstrauss transform. In *Proceedings of the forty-second ACM symposium on Theory of computing*, pages 341–350, 2010.
- 22 Sanjoy Dasgupta. Learning mixtures of gaussians. In *40th Annual Symposium on Foundations of Computer Science (Cat. No. 99CB37039)*, pages 634–644. IEEE, 1999.
- 23 Sanjoy Dasgupta and Anupam Gupta. An elementary proof of the Johnson-Lindenstrauss lemma. *International Computer Science Institute, Technical Report*, 22(1):1–5, 1999.
- 24 Victor H de la Peña and Stephen J Montgomery-Smith. Decoupling inequalities for the tail probabilities of multivariate u-statistics. *The Annals of Probability*, pages 806–816, 1995.
- 25 Devdatt P Dubhashi and Desh Ranjan. Balls and bins: A study in negative dependence. *BRICS Report Series*, 3(25), 1996.
- 26 Morris L Eaton. A note on symmetric bernoulli random variables. *The annals of mathematical statistics*, 41(4):1223–1226, 1970.
- 27 Cees M Fortuin, Pieter W Kasteleyn, and Jean Ginibre. Correlation inequalities on some partially ordered sets. *Communications in Mathematical Physics*, 22(2):89–103, 1971.
- 28 Peter Frankl and Hiroshi Maehara. The Johnson–Lindenstrauss lemma and the sphericity of some graphs. *Journal of Combinatorial Theory, Series B*, 44(3):355–362, 1988.
- 29 Peter Frankl and Hiroshi Maehara. Some geometric applications of the beta distribution. *Annals of the Institute of Statistical Mathematics*, 42(3):463–474, 1990.
- 30 Casper B Freksen, Lior Kamma, and Kasper Green Larsen. Fully understanding the hashing trick. In *Advances in Neural Information Processing Systems*, pages 5389–5399, 2018.
- 31 David Lee Hanson and Farroll Tim Wright. A bound on tail probabilities for quadratic forms in independent random variables. *The Annals of Mathematical Statistics*, 42(3):1079–1083, 1971.
- 32 G.H. Hardy, Karreman Mathematics Research Collection, J.E. Littlewood, G. Pólya, G. Pólya, and D.E. Littlewood. *Inequalities*. Cambridge Mathematical Library. Cambridge University Press, 1952. URL: <https://books.google.at/books?id=t1RCSP8YKt8C>.
- 33 Pawel Hitczenko. Domination inequality for martingale transforms of a rademacher sequence. *Israel Journal of Mathematics*, 84(1-2):161–178, 1993.
- 34 Wassily Hoeffding. Probability inequalities for sums of bounded random variables. In *The Collected Works of Wassily Hoeffding*, pages 409–426. Springer, 1994.

- 35 Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing*, pages 604–613, 1998.
- 36 Meena Jagadeesan. Understanding sparse jl for feature hashing. In *Advances in Neural Information Processing Systems*, pages 15203–15213, 2019. [arXiv:1903.03605](#).
- 37 Thathachar S Jayram and David P Woodruff. Optimal bounds for Johnson–Lindenstrauss transforms and streaming problems with subconstant error. *ACM Transactions on Algorithms (TALG)*, 9(3):1–17, 2013.
- 38 William B Johnson and Joram Lindenstrauss. Extensions of lipschitz mappings into a Hilbert space. *Contemporary mathematics*, 26(189-206):1, 1984.
- 39 William B Johnson and Assaf Naor. The Johnson–Lindenstrauss lemma almost characterizes Hilbert space, but not quite. *Discrete & Computational Geometry*, 43(3):542–553, 2010.
- 40 Daniel Kane, Raghu Meka, and Jelani Nelson. Almost optimal explicit Johnson–Lindenstrauss families. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*, pages 628–639. Springer, 2011.
- 41 Daniel M Kane and Jelani Nelson. Sparser Johnson–Lindenstrauss transforms. *Journal of the ACM (JACM)*, 61(1):1–23, 2014.
- 42 Krishnaram Kenthapadi, Aleksandra Korolova, Ilya Mironov, and Nina Mishra. Privacy via the Johnson–Lindenstrauss transform. *Journal of Privacy and Confidentiality*, 5(1):39–71, 2013.
- 43 Michael Kerber and Sharath Raghvendra. Approximation and streaming algorithms for projective clustering via random projections. *arXiv preprint*, 2014. [arXiv:1407.2063](#).
- 44 Aleksandr Khintchine. Über dyadische brüche. *Mathematische Zeitschrift*, 18(1):109–116, 1923.
- 45 Andreas Knoblauch. Closed-form expressions for the moments of the binomial probability distribution. *SIAM Journal on Applied Mathematics*, 69(1):197–204, 2008.
- 46 Konrad Kolesko and Rafał Latała. Moment estimates for chaoses generated by symmetric random variables with logarithmically convex tails. *Statistics & Probability Letters*, 107:210–214, 2015.
- 47 Samory Kpotufe and Bharath Sriperumbudur. Gaussian sketching yields a jl lemma in rkhs. In *International Conference on Artificial Intelligence and Statistics*, pages 3928–3937, 2020.
- 48 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images, 2009.
- 49 Rafał Latała. Tail and moment estimates for some types of chaos. *Studia mathematica*, 135(1):39–53, 1999.
- 50 Rafał Latała et al. Estimation of moments of sums of independent real random variables. *The Annals of Probability*, 25(3):1502–1513, 1997.
- 51 Ping Li, Trevor J Hastie, and Kenneth W Church. Very sparse random projections. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 287–296, 2006.
- 52 Nathan Linial, Eran London, and Yuri Rabinovich. The geometry of graphs and some of its algorithmic applications. *Combinatorica*, 15(2):215–245, 1995.
- 53 John E Littlewood. On the probability in the tail of a binomial distribution. *Advances in Applied Probability*, 1(1):43–72, 1969.
- 54 Konstantin Makarychev, Yury Makarychev, and Ilya 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*, pages 1027–1038, 2019.
- 55 Jiří Matoušek. On variants of the Johnson–Lindenstrauss lemma. *Random Structures & Algorithms*, 33(2):142–156, 2008.
- 56 Brendan D McKay. On littlewood’s estimate for the binomial distribution. *Advances in Applied Probability*, 21(2):475–478, 1989.

- 57 Jelani Nelson and Huy L. Nguyễn. Osnap: Faster numerical linear algebra algorithms via sparser subspace embeddings. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 117–126. IEEE, 2013.
- 58 Alfréd Rényi et al. On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*. The Regents of the University of California, 1961.
- 59 Herbert Robbins. A remark on Stirling’s formula. *The American mathematical monthly*, 62(1):26–29, 1955.
- 60 Mark Rudelson, Roman Vershynin, et al. Hanson-wright inequality and sub-gaussian concentration. *Electronic Communications in Probability*, 18, 2013.
- 61 Tamas Sarlos. Improved approximation algorithms for large matrices via random projections. In *2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS’06)*, pages 143–152. IEEE, 2006.
- 62 Pantelimon Stanica. Good lower and upper bounds on binomial coefficients. *Journal of Inequalities in Pure and Applied Mathematics*, 2(3):30, 2001.
- 63 Roman Vershynin. A simple decoupling inequality in probability theory. *preprint*, 2011.
- 64 Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- 65 Kilian Weinberger, Anirban Dasgupta, John Langford, Alex Smola, and Josh Attenberg. Feature hashing for large scale multitask learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 1113–1120, 2009.
- 66 Alfred Witkowski. A new proof of the monotonicity of power means. *J. Ineq. Pure and Appl. Math*, 5(1), 2004.
- 67 Shuheng Zhou. Sparse Hanson–Wright inequalities for subgaussian quadratic forms. *Bernoulli*, 25(3):1603–1639, 2019. appears in 2015 at [arXiv:1510.05517](https://arxiv.org/abs/1510.05517).

A

 Some remarks on prior works

A.1 Some issues with numeric constants

Lemma 2.1 in [36] gives the following bound (expressed in our notation)

$$\|Z_r\|_d \lesssim \begin{cases} dp & d = 2 \text{ or } d \leq pe/v^2 \\ \min\left(\frac{d^2 v^2}{\log(dv^2/p)}, \frac{d}{\log(1/p)}\right) & 1 \leq \log(dv^2/p) \leq d \\ v^2 (p/dv^2)^{2/d} & d < \log(dv^2/p) \end{cases}$$

There is a minor mistake in splitting the branches: they emerge from taking the derivative test of the function $d^2 v^2 u^{-2} (p/dv^2)^{1/u}$ where $1 \leq u \leq d/2$ (Lemma D.1). Here the local maxima occurs at $u = \log(dv^2/p)/2$ and when comparing this with edges $u = 1$ and $u = d/2$ we obtain the conditions $2 \leq \log(dv^2/p)$ and $\log(dv^2/p) \leq d$. Thus, the splitting conditions should be a bit different; this particular issue doesn’t affect the bounds expressed in the asymptotic notation; we report it with intent to motivate our effort in giving a simple and clear proof.

A.2 Gaps in symmetrization

Section 2.2 of [36], when explaining the proof strategy, proposes to apply the bounds on Z_r defined in Equation (10) assuming they are symmetric. But Z_r are not symmetric (it is easy to see they have positive higher-order moments), thus extra work is needed to push this argument forward.

B Concluding Main Theorem

Without losing generality, we assume that $d = \log(1/\delta)$ is even. Recall that we denote $v = v_d(x)$, also without losing generality we assume that v^{-2} is an integer. For $K = v^{-2}$ define the following quantities

$$\begin{aligned} I_1 &\triangleq \max_q \left\{ d/q \cdot (m/d)^{1/q} \cdot qp : \log(q/Kp) \leq q/K \leq 2, 2 \leq q \leq d \right\} \\ I_2 &\triangleq \max_q \left\{ d/q \cdot (m/d)^{1/q} \cdot K(Kp^{2K/q})^2 : \log(q/Kp) \leq 2 \leq q/K, 2 \leq q \leq d \right\} \\ I_3 &\triangleq \max_q \left\{ d/q \cdot (m/d)^{1/q} \cdot K^{-1}q^2 / \log^2(q/Kp) : \max(2, q/K) \leq \log(q/Kp) \leq q, 2 \leq q \leq d \right\} \\ I_4 &\triangleq \max_q \left\{ d/q \cdot (m/d)^{1/q} \cdot K^{-1}(Kp)^{2/q} : q \leq \log(q/Kp), 2 \leq q \leq d \right\}. \end{aligned}$$

Following the proof outline we arrive at Corollary 13. Taking into account Lemma 11 and Lemma 9, implies that:

$$\left\| \sum_{r=1}^m Z_r \right\|_d \leq O(\max(I_1, I_2, I_3, I_4)).$$

The goal is to prove that for $t = s\epsilon$ we have

$$\left\| \sum_{r=1}^m Z_r \right\|_d \leq t/e, \tag{16}$$

and then the result follows from Markov's inequality. We give first bounds for I_1, I_2, I_4 as they are fairly easy to obtain. The case of I_3 is analyzed as the last one.

B.1 First Branch

We will show the following bound

► **Lemma 20.** *We have*

$$I_1 \leq O(dmp^2)^{1/2}.$$

Proof of Lemma 20. We have

$$\begin{aligned} I_1 &= \max_q \left\{ pd(m/d)^{1/q} : \log(q/Kp) \leq q/K \leq 2, 2 \leq q \leq d \right\} \\ &\leq (dmp^2)^{1/2} \end{aligned}$$

where the inequality follows because $m \geq d$ and $1/q \leq \frac{1}{2}$ (for q satisfying the constraints). This completes the proof. ◀

B.2 Second Branch

We will show the following bound

► **Lemma 21.** *For $p \leq 2e^{-2}$ we have*

$$I_2 \leq (dmp^2)^{1/2}.$$

18:16 Sparse Random Projections

Proof of Lemma 20. For q satisfying the constraint we have $K/q \geq e^{-2}/p$ which, due to $p \leq 2e^{-2}$, implies $K/q \geq 1/2$. Then $p^{2K/q} \leq p$ (recall that $p < 1!$) and thus

$$I_2 \leq \max_q \left\{ d/q \cdot (m/d)^{1/q} \cdot Kp : \log(q/Kp) \leq 2 \leq q/K, 2 \leq q \leq d \right\}.$$

For q within the constraints we have $K/q \leq \frac{1}{2}$ and therefore

$$I_2 \leq \frac{p}{2} \max_q \left\{ d \cdot (m/d)^{1/q} : \log(q/Kp) \leq 2 \leq q/K, 2 \leq q \leq d \right\}.$$

Since $m/d \geq 1$ the expression under the maximum decreases with q , thus is not bigger than the value at $q = 2$. Thus, $I_2 \leq p(dm)^{1/2}/2$ and the result follows. \blacktriangleleft

B.3 Fourth Branch

We will prove the following bound

► **Lemma 22.** *We have*

$$I_4 \leq \begin{cases} (dmp^2)^{1/2} & \log(dv^4/mp^2) \leq 2 \\ dv^2/\log(dv^4/mp^2) & \log(dv^4/mp^2) > 2 \end{cases}.$$

Proof of Lemma 22. We have

$$I_4 = \max_q \left\{ K^{-1} \cdot d/q \cdot (K^2 p^2 m/d)^{1/q} : q \leq \log(q/Kp), 2 \leq q \leq d \right\}.$$

Let $a = K^2 p^2 m/d$, the expression under the maximum is proportional to $1/q \cdot a^{1/q}$. We now apply Proposition 15: for $a \geq 1$ the maximum is not bigger than the value at $q = 2$, so

$$I_4 \leq (dmp^2)^{1/2}.$$

We now can assume $a < 1$, equivalent to $K^2 p^2 m < d$. The global maximum is at $q = \log(1/a)$, thus our maximum is still at $q = 2$ when $\log(1/a) \leq 2$ and otherwise is not bigger than the value at $q = \log(1/a)$. We then obtain

$$I_4 \leq K^{-1} d / \log(d/mp^2 K^2) \leq K^{-1} d = dv^2.$$

This complete the proof. \blacktriangleleft

B.4 Third Branch

We will show the following bound

► **Lemma 23.** *Suppose that $v^2 \geq \epsilon/d^2$, then*

$$I_3 \leq O(dmp^2)^{1/2} + O(dv/\log(dv^2/p))^2$$

Proof of Lemma 23. The proof is based on splitting the maximum into three regimes: $q \in [2, 3], 3 \leq q \leq \log(m/d)$ and $\log(m/d) \leq q \leq d$. Define

$$I^0 = \max_q \left\{ d/q \cdot (m/d)^{1/q} \cdot v^2 q^2 / \log^2(qv^2/p) : 2 \leq \log(qv^2/p) \leq q \leq d, 2 \leq q \leq 3 \right\}$$

$$I^- = \max_q \left\{ d/q \cdot (m/d)^{1/q} \cdot v^2 q^2 / \log^2(qv^2/p) : 2 \leq \log(qv^2/p) \leq q \leq d, 3 \leq q \leq \log(m/d) \right\}$$

$$I^+ = \max_q \left\{ d/q \cdot (m/d)^{1/q} \cdot v^2 q^2 / \log^2(qv^2/p) : 2 \leq \log(qv^2/p) \leq q \leq d, \log(m/d) \leq q \leq d \right\}.$$

so that we have $I_3 \leq \max(I^0, I^+, I^-)$ (for convenience, we replace the constraint $\max(2, qv^2) \leq \log(qv^2/p)$ in I_3 by the weaker one $2 \leq \log(qv^2/p)$). By the assumptions we have $v^2/p \geq m\epsilon/d^2$. Since $m \geq d\epsilon^{-2}$ we have $\epsilon \geq (d/m)^{1/2}$, and thus

$$v^2/p \geq (m/d)^{1/2} \cdot d^{-1}.$$

▷ **Claim 24.** We have $I^- \leq O(d^2v^2/\log^2(dv^2/p))$ when $\log d \leq \frac{5\log(m/d)}{12}$.

Proof of Claim. For any q satisfying the restrictions it holds that

$$\begin{aligned} q &\geq \log(v^2/p) \\ &\geq \frac{\log(m/d)}{2} - \log d \\ &\geq \frac{\log(m/d)}{12}. \end{aligned}$$

We then have $(m/d)^{1/q} \leq O(1)$ and thus

$$I^- \leq \max_q \{d \cdot qv^2/\log^2(qv^2/p) : 2 \leq \log(qv^2/p) \leq q \leq d, 3 \leq q \leq \log(m/d)\}.$$

Considering the auxiliary function $u \rightarrow u/\log^2 u$ with $u = qv^2/p \geq e^2$, we see that it decreases in u and hence in q for fixed v^2 and p . The expression is thus not smaller than its value at $q = d$, which gives

$$I^- \leq d^2v^2/\log^2(dv^2/p),$$

and completes the proof. ◁

▷ **Claim 25.** We have $I^- \leq d^2v^2/\log^2(dv^2/p)$ when $\log d > \frac{5\log(m/d)}{12}$.

Proof of Claim. We have that $dv^2/p \geq m\epsilon/d \geq (m/d)^{1/2}$ and therefore

$$\begin{aligned} I^- &\leq dv^2d(m/d)^{1/3}\log(m/d) \\ &\leq dv^2(m/d)^{5/12}/\log^2(m/d) \\ &\leq dv^2(m/d)^{5/12}/\log^2(dv^2/p) \\ &\leq O(d^2v^2/\log^2(dv^2/p)), \end{aligned}$$

which completes the proof. ◁

▷ **Claim 26.** We have $I^+ \leq O(d^2v^2/\log^2(dv^2/p))$

Proof of Claim. We have $(m/d)^{1/q} \leq e$ for $q \geq \log(m/d)$, thus

$$I^+ \leq d \cdot \max_q \{qv^2/\log^2(qv^2/p) : 2 \leq \max(\log(qv^2/p), \log(m/d)) \leq q \leq d\}.$$

Considering the auxiliary function $u \rightarrow u/\log^2 u$ with $u = qv^2/p \geq e^2$, we see that it decreases in u and hence in q for fixed v^2 and p . The expression is thus not smaller than its value at $q = d$, which gives

$$I^+ \leq O(d^2v^2/\log^2(dv^2/p))$$

and the claim follows. ◁

18:18 Sparse Random Projections

▷ **Claim 27.** We have $I^0 \leq O((dmp^2)^{1/2})$.

Proof of Claim. We have $I^0 \leq O(v^2(md)^{1/2})$ because $(m/d)^{1/q} \leq (m/d)^{1/2}$ (due to $m/d \geq 1$ and $q \geq 2$). However, for $q \in [2, 3]$ the constraint $\log(qv^2/p) \leq q$ gives $v^2 \leq O(p)$. Thus

$$I^0 \leq O(p(md)^{1/2}),$$

which completes the proof. ◁

The result follows now by combining the above three claims. ◀

B.5 Merging Branch Bounds

To conclude the main result it suffices to satisfy

$$c \cdot \max(I_1, I_2, I_3, I_4) \leq s\epsilon \tag{17}$$

for some absolute constant c . The condition in Equation (17) for I_1, I_2 is equivalent to $c \cdot (dmp^2)^{1/2} \leq s\epsilon$, which holds when

$$m \geq \Omega(d\epsilon^{-2}). \tag{18}$$

To satisfy Equation (17) for I_4 we require, in addition to Equation (18), that $cdv^2 \leq s\epsilon$, equivalent to

$$v \leq O((s\epsilon)^{1/2}/d^{1/2}). \tag{19}$$

Finally, in order to satisfy Equation (17) for I_3 we observe that, under the restriction

$$v^2 \geq s\epsilon/d^2, \tag{20}$$

the bound in Lemma 23 gives

$$I_3 \leq O(dmp^2)^{1/2} + O(dv/\log(m\epsilon/d))^2,$$

which follows because $\log(dv^2/p) \geq \log(s\epsilon/dp) = \log(m\epsilon/d)$. Thus, in addition to Equation (18) and Equation (20) it suffices that

$$v \leq O((s\epsilon)^{1/2} \log(m\epsilon/d)/d). \tag{21}$$

Now observe that for

$$v = \Theta(s\epsilon)^{1/2} \min(\log(m\epsilon/d)/d, 1/d^{1/2}) \tag{22}$$

the condition in Equation (20) is automatically satisfied. Thus, the theorem holds for v as above, and clearly for any smaller v .