

1 High Probability Tensor Sketch

2 Thomas Dybdahl Ahle

3 IT University of Copenhagen

4 Jakob Bæk Tejs Knudsen

5 Copenhagen University

6 — Abstract —

7 We construct a structured Johnson Lindenstrauss transformation that can be applied to simple
8 tensors on the form $x = x^{(1)} \otimes \dots \otimes x^{(c)} \in \mathbb{R}^{d^c}$ in time nearly cd . That is, exponentially faster than
9 writing out the Kronecker product and then mapping down.

10 These matrices, M , which preserves the norm of any $x \in \mathbb{R}^{d^c}$, such that $\| \|Mx\|_2 - \|x\|_2 \| \leq \epsilon$
11 with probability $1 - \delta$, can be taken to have just $\tilde{O}(c^2 \epsilon^{-2} (\log 1/\delta)^3)$ rows. This is within $c^2 (\log 1/\delta)^2$
12 of optimal for any JL matrix [16], and improves upon earlier ‘Tensor Sketch’ constructions by Pagh
13 and Pham [24, 25], which used $\tilde{O}(3^c \epsilon^{-2} \delta^{-1})$ rows, by an exponential amount in both c and δ^{-1} .

14 It was shown by Avron, Nguyen and Woodruff in [5] that Tensor Sketch is a subspace embedding.
15 This has a large number of applications [28], such as guaranteeing the correctness of kernel-linear
16 regression performed directly on the reduced vectors. We show that our construction is a subspace
17 embedding too, improving again upon the exponential dependency on c and δ^{-1} , enabling sketching
18 of much higher order polynomial kernels, such as Taylor approximations to the ubiquitous Gaussian
19 radial basis function.

20 Technically, we construct our matrix M such that $M(x \otimes y) = Tx \circ T'y$ where \circ is the Hadamard
21 (element-wise) product and T and T' support fast matrix-vector multiplication ala [1]. To analyze
22 the behavior of Mx on non-simple x , we show a higher order version of Khintchine’s inequality [14],
23 related to the higher order Gaussian chaos analysis by Latała [19, 17]. Finally we show that such
24 sketches can be combined recursively, in a way that doesn’t increase the dependency on c by much.

25 **2012 ACM Subject Classification** Theory of computation → Sketching and sampling

26 **Keywords and phrases** Tensors, Sketching, Johnson Lindenstrauss

27 **Digital Object Identifier** 10.4230/LIPIcs...

28 **Acknowledgements** Anders Aamand

29 **1** Introduction

30 The polynomial method has recently found many great applications in algorithm design, such
31 as finding orthogonal vectors [2] and gap amplification in nearest neighbour [27]. Consider
32 the polynomial $P(x, y) = \sum_{i < j < k} (x_i + y_i - x_i y_i)(x_j + y_j - x_j y_j)(x_k + y_k - x_k y_k)$ which
33 counts the number of triangles in the union between two graphs, $x, y \in \{0, 1\}^{\binom{n}{2}}$, expressed
34 as binary vectors over the edges. Splitting P into monomials, one may construct functions
35 $f, g : \{0, 1\}^{\binom{n}{2}} \rightarrow \mathbb{R}^m$ such that $P(x, y) = \langle f(x), g(y) \rangle$, and use these embeddings to solve
36 the ‘most triangles in union’ problem on a database of graphs, using an off the shelf nearest
37 neighbours algorithm. (See our Applications for more on this.)

38 Other examples are kernel functions in statistics, such as $P(x, y) = \exp(-\|x - y\|_2^2) =$
39 $\sum_{k \geq 0} (-1)^k \|x - y\|_2^{2k} / k!$, the Gaussian Radial Basis Function. The celebrated ‘kernel trick’
40 has been used in linear methods such as kernel PCA kernel nearest neighbour or kernel
41 regression to allow detection of nonlinear dependencies between data without explicitly
42 constructing feature vectors in high dimensional spaces. However the ‘trick’ requires the
43 computation of the inner product between all pairs of points, while explicit embeddings scale
44 linearly in the number of points, and have thus recently experienced a comeback.



© Thomas Dybdahl Ahle, Jakob Bæk Tejs Knudsen;
licensed under Creative Commons License CC-BY

Leibniz International Proceedings in Informatics

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

45 While these polynomial expansions often produce prohibitively large vectors, they can
 46 often be reduced by some means, such as the Johnson-Lindenstrauss transform [13]. This
 47 in turn creates a strange phenomenon where we first blow up the dimension only to later
 48 squash it back down. It is tempting to look for shortcut to go straight to the final dimension.

49 This was the idea of Pagh and Pham [24, 25] with Tensor Sketch. The observation was
 50 that $P(x, y) = \langle x, y \rangle^2 = \langle x \otimes x, y \otimes y \rangle$, where $x \otimes x$ is the tensor product (Kronecker) of x
 51 with itself. They further observed that if C and C' are independent Count Sketch matrices,
 52 then $\langle Cx * C'x, Cy * C'y \rangle \approx \langle x, y \rangle^2$ while the dimension of $Cx * C'x$ is much smaller than
 53 of $x \otimes x$. Since the convolution $(\cdot * \cdot)$ can be computed in near linear time using the fast
 54 Fourier transformation, they could sketch $x \otimes x$ in basically the time required to sketch x .
 55 For a higher order polynomial kernel, replacing $f(x) = x^{\otimes c}$ with $f(x) = C^{(1)}x * \dots * C^{(c)}x$
 56 thus takes the sketching time from d^c to cd . A huge improvement!

57 In this paper we improve on the main shortcomings of Tensor Sketch: To preserve the
 58 norm of vectors up to $1 \pm \epsilon$ with probability $1 - \delta$, it requires embedding into dimension
 59 roughly $3^c \epsilon^{-2} \delta^{-1}$. The exponential dependency on c greatly limits the degree of polynomials
 60 that can be embedded, and the linear dependency on δ^{-1} means we can't use a standard
 61 union bound trick to get e.g. a near neighbour preserving embedding [11], as could be
 62 achieved with the Johnson Lindenstrauss transform, which embeds into only $\epsilon^{-2} \log 1/\delta$
 63 dimensions. We overcome both of these obstacles, by analyzing a scheme, that with the same
 64 embedding time, requires only $c^2 \epsilon^{-2} (\log 1/\delta)^3$ dimensions.

65 A hugely important idea was introduced by Avron et al. [5]: They proved that a Tensor
 66 Sketch with sufficiently many rows is a Subspace embedding. This allowed many applica-
 67 tion that previously were only applied heuristically, such as solving a regression problem
 68 $\arg \min_x \|Ax - y\|_2$ directly in the reduced space while guaranteeing correct results. Using
 69 the subspace embeddings, they obtained the fastest known algorithms for computing an
 70 computing an approximate kernel PCA and many other problems.

71 However, the weaknesses of Tensor Sketch remained: The exponential dependency on c
 72 meant that the method could only be applied with relatively low degree polynomials. In this
 73 paper we show that our High Probability Tensor Sketch is also a subspace embedding, solving
 74 this major roadblock. We also show a second version of our sketch, which improves upon [5]
 75 by allowing an embedding dimension linear in the subspace dimension, rather than quadratic.
 76 In many uses of the subspace method the embedding dimension becomes larger than the
 77 number of points, which means we can get a quadratic improvement on these applications.

78 Our approach is to analyze fast family of Johnson Lindenstrauss matrices M with the
 79 further property that $M(x \otimes y) = M'x \circ M''y$ where \circ is the Hadamard (or element-wise)
 80 vector product. We also analyze the case where M' and M'' are fully independent Gaussian
 81 matrices, and show that we are within a single factor $\log 1/\delta$ in embedding dimension
 82 while supporting much faster matrix-vector multiplication. The direct application of this
 83 method, $M^{(1)}x^{(1)} \circ \dots \circ M^{(c)}$, would result in an exponential dependency on c , but by instead
 84 combining vectors as $M^{(1)}x^{(1)} \circ M^{(1)}(M^{(2)}x^{(2)}) \circ M^{(2)}(\dots$ we show that this dependency
 85 can be reduced to c^2 . See also the Technical Overview below.

86 1.1 Overview

87 Our main contribution is to answer the questions “Does Tensor Sketch work with high
 88 probability?” and “Does there exist subspace embeddings for higher order polynomial
 89 embeddings?” For both of those questions, the answer is yes!

90 ► **Theorem 1 (Construction A).** *There is a distribution \mathcal{M} over matrices $M \in \mathbb{R}^{m \times d^c}$ where*

- 91 1. $\|Mx\|_2 - \|x\|_2 \leq \epsilon$ with probability $\geq 1 - \delta$ for any $x \in \mathbb{R}^{d^c}$.
- 92 2. M can be applied to tensors $x^{(1)} \otimes \dots \otimes x^{(c)} \in \mathbb{R}^{d^c}$ in time $O(c(d \log d + m \log m))$.
- 93 3. m can be taken to be $O(c^2 \epsilon^{-2} (\log 1/\delta) (\log 1/\epsilon \delta)^2)$.
- 94 4. \mathcal{M} can compute fast approximate matrix multiplication: $\Pr[\|A^T M^T M B - A^T B\|_F >$
95 $\epsilon \|A\|_F \|B\|_F] < \delta$.
- 96 5. There is an $m = O(c^2 \epsilon^{-2} \lambda^2 (\log 1/\delta) (\log 1/\epsilon \delta)^2)$ such that \mathcal{M} is an (ϵ, δ) -subspace embed-
97 ding. (See definition 11.)

98 The result matches, up to a single factor $\log 1/\delta$ the embedding dimension needed for
99 fully independent Gaussian matrices M and M' , for $\|Mx \circ M'y\|_2$ to approximate $\|x \otimes y\|_2$.
100 (See Appendix, Theorem 27.) However, suffering slightly in the embedding time, we can go
101 all the way down to one:

102 ► **Theorem 2** (Construction B). *There is a distribution \mathcal{M} over matrices $M \in \mathbb{R}^{m \times d^c}$ where*

- 103 1. $\|Mx\|_2 - \|x\|_2 \leq \epsilon$ with probability $\geq 1 - \delta$ for any $x \in \mathbb{R}^{d^c}$.
- 104 2. Matrices $M \sim \mathcal{M}$ can be applied to tensors $x^{(1)} \otimes \dots \otimes x^{(c)} \in \mathbb{R}^{d^c}$ in time $O(cm \min(d, m))$.
- 105 3. m can be taken to be $O(c^2 \epsilon^{-2} \log 1/\delta \log^2(c \epsilon^{-1} \log 1/\delta))$.
- 106 4. There is an $m = O(c^2 \epsilon^{-2} (\lambda + \log 1/\delta) \log^2(c \epsilon^{-1} \lambda \log 1/\delta))$ such that \mathcal{M} is an (ϵ, δ) -
107 subspace embedding. (See definition 11.)

108 While the first theorem requires an intricate analysis of the combination of two Fast-JL
109 matrices, the second one follows nearly directly from our general recursive construction
110 theorem below:

111 ► **Theorem 3.** *Let $c > 0$ be a positive integer, and $Q^{(1)} \in \mathbb{R}^{m \times d}$ and $Q^{(i)} \in \mathbb{R}^{m \times md}$ be
112 independent random matrices for every $i \in [c] \setminus \{1\}$. Define $M^{(k)} = Q^{(k)}(M^{(k-1)} \otimes I_d) \in$
113 $\mathbb{R}^{m \times d^k}$ for $k \in [c]$, where $M^{(0)} = 1 \in \mathbb{R}$. Let $t > 0$ be a positive integer, and let $k_i \in [c]$ for
114 every $i \in [t]$. Then the matrix*

$$115 \quad M = \bigoplus_{i \in [t]} M^{(k_i)} \in \mathbb{R}^{tm \times \sum_{i \in [t]} d^{k_i}}$$

116 *has the following properties.*

- 117 1. Let $\epsilon \in (0, 1)$ and $\delta > 0$. If $Q^{(i)}$ has $(\epsilon/2c, \delta/c)$ -JL property for every $i \in [c]$, then M has
118 (ϵ, δ) -JL property.
- 119 2. If $Q^{(i)}x$ can be evaluated in time T , where $x \in \mathbb{R}^{md}$, for every $i \in [c] \setminus \{1\}$, and $Q^{(1)}y$
120 can be evaluated in time T' , where $y \in \mathbb{R}^d$, then $M(\bigoplus_{i \in [t]} \bigotimes_{j \in [k_i]} x^{(i,j)})$ can be evaluated
121 in time $O(T't + T \sum_{i \in [t]} k_i)$, where $x^{(i,j)} \in \mathbb{R}^d$ for every $i \in [t], j \in [k_i]$.

122 Now the difference between Construction A and Construction B is simply which matrices
123 $Q^{(1)}, \dots, Q^{(c)}$ that are used as basis for the construction.

124 Paper Structure

125 The paper is structured as follows: After the comparison to related work and preliminaries
126 we give a Technical overview of the sketch. We find it is useful to have some established
127 notation before this section.

128 The technical part is split in three: We first show Theorem 3. This gives a recursive
129 construction, which can be applied to tensors using any of the shelf Johnson-Lindenstrauss
130 matrix. Combined with the fastest analysis of Fast-JL [15] this gives our theorem 2.

131 We proceed to analyze a small change in the construction of Fast-JL matrices, which
 132 allow for very fast application to tensor products. Specifically we show that the random
 133 diagonal matrix can be replaced by the Kronecker product of two smaller diagonal matrices
 134 without losing the JL-property, if the number of rows is increased slightly. This gives our
 135 theorem 1.

136 Finally in the last section, we show some algorithmic applications of our constructions.
 137 For example how to use it to find the two graphs in a database whose union has the most
 138 triangles.

139 **1.2 Related work**

140 Work related to sketching of tensors and explicit kernel embeddings is found in fields ranging
 141 from pure mathematics to physics and machine learning. Hence we only try to compare
 142 ourselves with the four most common types we have found.

143 We focus particularly on the work on subspace embeddings [25, 5], since it is most directly
 144 comparable to ours. An extra entry in this category is [4], which is currently in review, and
 145 which we were made aware of while writing this paper. That work is in double blind review,
 146 but by the time of the final version of this paper, we should be able to cite it properly.

147 **Subspace embeddings**

148 For most applications [5], the subspace dimension, λ , will be much larger than the input
 149 dimension, d , but smaller than the implicit dimension d^c . Hence the size of the sketch, m ,
 150 will also be assumed to satisfy $d \ll m \ll d^c$ for the purposes of stating the results. We will
 151 hide constant factors, and $\log 1/\epsilon$, $\log d$, $\log m$, $\log c$, $\log \lambda$ factors.

152 Note that we can always go from m down to $\approx \epsilon^{-2}(\lambda + \log 1/\delta)$ by applying a fast-JL
 153 transformation after embedding. This works because the product of two subspace embeddings
 154 is also a subspace embedding, and because fast-JL is a subspace embedding by the net-
 155 argument (see lemma 13). The embedding dimensions in the table should thus mainly be
 156 seen as a space dependency, rather than the actual embedding dimension for applications.

Reference	Embedding dimension, m	Embedding time	Note
[25, 5]	$\tilde{O}(3^c d \lambda^2 \delta^{-1} \epsilon^{-2})$	$\tilde{O}(c(d+m))$	
Theorem 1	$\tilde{O}(c^2 \lambda^2 (\log 1/\delta)^3 \epsilon^{-2})$	$\tilde{O}(c(d+m))$	
Theorem 2	$\tilde{O}(c^2 (\lambda + \log 1/\delta) \epsilon^{-2})$	$\tilde{O}(c d m)$	
[4], Theorem 1	$\tilde{O}(c \lambda^2 \delta^{-1} \epsilon^{-2})$	$\tilde{O}(c(d+m))$	Independent work.
[4], Theorem 2	$\tilde{O}(c^6 \lambda (\log 1/\delta)^5 \epsilon^{-2})$	$\tilde{O}(c(d+m))$	Independent work

158 Some of the results, in particular [25, 5] and [4] Theorem 1 can be applied faster when
 159 the input is sparse. Our results, as well as [4], Theorem 2 can similarly be optimized for
 160 sparse inputs, by preprocessing vectors with an implementation of Sparse JL [7].

161 In comparison to the previous result [25, 5] we are clearly better with an exponential
 162 improvement in c as well as δ .

163 Compared to the new work of [4], all four bounds have some region of superiority. Their
 164 first bound of has the best dependency on c , but has an exponential dependency on $\log 1/\delta$.
 165 Their second bound has an only linear dependency on $d + \lambda$, but has large polynomial
 166 dependencies on c and $\log 1/\delta$.

167 Technically the methods of all five bounds are similar, but some details and much of the
 168 analysis differ. Our results as well as the results of [4] use recursive constructions to avoid
 169 exponential dependency on c , however the shape of the recursion differs. We show all of

170 our results using the p -moment method, while [4] Theorem 1 and [25, 5] are shown using
 171 2nd-moment analysis. This explains much of why their dependency on δ is worse.

172 **Approximate Kernel Expansions**

173 A classic result by Rahimi and Recht [26] shows how to compute an embedding for any
 174 shift-invariant kernel function $k(\|x - y\|_2)$ in time $O(dm)$. In [18] this is improved to any
 175 kernel on the form $k(\langle x, y \rangle)$ and time $O((m + d) \log d)$. This is basically optimal in terms of
 176 time and space, however the method does not handle kernel functions that can't be specified
 177 as a function of the inner product, and it doesn't provide subspace embeddings. See also [22]
 178 for more approaches along the same line.

179 **Tensor Sparsification**

180 There is also a literature of tensor sparsification based on sampling [23], however unless
 181 the vectors tensored are already very smooth (such as ± 1 vectors), the sampling has to be
 182 weighted by the data. This means that these methods in aren't applicable in general to the
 183 types of problems we consider, where the tensor usually isn't known when the sketching
 184 function is sampled.

185 **Hyper-plane rounding**

186 An alternative approach is to use hyper-plane rounding to get vectors on the form ± 1 . Let
 187 $\rho = \frac{\langle x, y \rangle}{\|x\| \|y\|}$, then we have $\langle \text{sign}(Mx), \text{sign}(My) \rangle = \sum_i \text{sign}(M_i x) \text{sign}(M_i y) = \sum_i X_i$, where
 188 X_i are independent Rademachers with $\mu/m = E[X_i] = 1 - \frac{2}{\pi} \arccos \rho = \frac{2}{\pi} \rho + O(\rho^3)$. By
 189 tail bounds then $\Pr[|\langle \text{sign}(Mx), \text{sign}(My) \rangle - \mu| > \epsilon \mu] \leq 2 \exp(-\min(\frac{\epsilon^2 \mu^2}{2\sigma^2}, \frac{3\epsilon \mu}{2}))$. Taking
 190 $m = O(\rho^{-2} \epsilon^{-2} \log 1/\delta)$ then suffices with high probability. After this we can simply sample
 191 from the tensor product using simple sample bounds.

192 The sign-sketch was first brought into the field of data-analysis by [6] and [27] was the
 193 first, in our knowledge, to use it with tensoring. The main issue with this approach is that
 194 it isn't a linear sketch, which hinders some applications, like subspace embeddings. It also
 195 takes dm time to calculate Mx and My . In general we would like fast-matrix-multiplication
 196 type results.

197 **2 Preliminaries**

198 We will use the following notation

- k, i, j Indices
- c Tensor order
- d Original dimension
(Assumed to be a power of 2.)
- d^c Implicit dimension
- m Sketch dimension
- λ Subspace dimension
- Λ Subspace of R^{d^c}
- M Sketching matrix
- \mathcal{M} Distribution of sketching matrices

We say $f(x) \lesssim g(x)$ if $f(x) = O(g(x))$.
 For $p \geq 1$ and random variables $X \in R$,
 we write $\|X\|_p = (E|X|^p)^{1/p}$. Note that
 $\|X + Y\|_p \leq \|X\|_p + \|Y\|_p$ by the Minkowski
 Inequality.

► **Definition 4** (Direct sum). *We define the direct sum of two vectors as*

$$x \oplus y = \begin{bmatrix} x \\ y \end{bmatrix},$$



and the direct sum between two matrices as

$$A \oplus B = \begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix}.$$

► **Definition 5** (Kronecker (tensor) product). We define the tensor-product (or Kronecker) of two matrices as:

$$A \otimes B = \begin{bmatrix} a_{1,1}B & \cdots & a_{1,n}B \\ \vdots & \ddots & \vdots \\ a_{m,1}B & \cdots & a_{m,n}B \end{bmatrix},$$

and in particular of two vectors: $x \otimes y = [x_1y_1, x_1y_2, \dots, x_ny_n]^T$. Taking the tensor-product of a vector with itself, we get the tensor-powers:

$$x^{\otimes k} = \underbrace{x \otimes \cdots \otimes x}_{k \text{ times}}$$

The Kronecker product has the useful mixed-product property when the sizes match up:

$$(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$$

We note in particular the vector variants $(I \otimes B)(x \otimes y) = x \otimes By$ and $\langle x \otimes y, z \otimes t \rangle = \langle x, y \rangle \langle z, t \rangle$.

► **Definition 6** (Hadamard product). Also sometimes known as the ‘element-wise product’:

$$x \circ y = [x_1y_1, x_2y_2, \dots, x_ny_n]^T.$$

Taking the Hadamard product with itself gives the Hadamard-power:

$$x^{\circ k} = \underbrace{x \circ \cdots \circ x}_{k \text{ times}} = [x_1^k, x_2^k, \dots, x_n^k]^T.$$

Definitions

► **Definition 7** (JL-moment property). We say a distribution over random matrices $M \in \mathbb{R}^{m \times d}$ has the (ϵ, δ) -JL-moment property, when

$$\| \|Mx\|_2^2 - 1 \|_p \leq \epsilon \delta^{1/p}$$

for all $p > 1$ and $x \in \mathbb{R}^d$, $\|x\| = 1$.

Note that by Markov’s inequality, the JL-moment-property implies $E\|Mx\|_2 = \|x\|_2$ and that taking $m = O(\epsilon^{-2} \log 1/\delta)$ suffices to have $\Pr[\|Mx\|_2 - \|x\|_2 > \epsilon] < \delta$ for any $x \in \mathbb{R}^d$. (This is sometimes known as the Distributional-JL property.)

► **Definition 8** (ϵ, δ) -Approximate Matrix Multiplication). We say a distribution over random matrices $M \in \mathbb{R}^{k \times d}$ has the (ϵ, δ) -Approximate Matrix Multiplication property if for any matrices A, B with proper dimensions,

$$\begin{aligned} & \| \|A^T M^T M B - A^T B \|_F \|_p \\ & \leq \epsilon \delta^{1/p} \|A\|_F \|B\|_F. \end{aligned}$$

► **Lemma 9** (Shown in [28]). Any distribution that has the (ϵ, δ) -JL-moment-property has the $(3\epsilon, \delta)$ -Approximate Matrix Multiplication property.

We note that the factor of 3 on ϵ can be removed by combining the analysis in [28] with Appendix Lemma 30.

► **Definition 10** (ϵ) -Subspace embedding). $M \in \mathbb{R}^{k \times D}$ is a subspace embedding for $\Lambda \subseteq \mathbb{R}^D$ if for any $x \in \Lambda$,

$$\| \|Mx\|_2 - \|x\|_2 \| \leq \epsilon.$$

► **Definition 11** (λ, ϵ) -Oblivious Subspace Embedding). A distribution, \mathcal{M} , over $\mathbb{R}^{m \times D}$ matrices is a (D, λ) -Oblivious Subspace Embedding if for any linear subspace, $\Lambda \subseteq \mathbb{R}^D$, of dimension λ , $M \sim \mathcal{M}$ is an ϵ -subspace embedding for Λ with probability at least $1 - \delta$.

► **Lemma 12**. Any distribution that has the $(\epsilon/(3\lambda), \delta)$ -JL-moment-property is a (λ, ϵ) -oblivious subspace embedding.

Proof. Let $U \in \mathbb{R}^{\lambda \times m}$ be orthonormal such that $U^T U = I$, it then suffices (by [28]) to show $\| \|U^T M^T M U - I \| \leq \epsilon$.

From lemma 9 we have that $\| \|U^T M^T M U - I \| \leq 3\epsilon \delta^{1/p} \|U\|_F^2 = 3\epsilon \delta^{1/p} \lambda$. ◀

► **Lemma 13**. There is a $C > 0$, such that any distribution that has the $(\epsilon, \delta e^{C\lambda})$ -JL-moment-property is a (λ, ϵ) -oblivious subspace embedding.

Proof. For any λ -dimensional subspace, Λ , there exists an ϵ -net $T \subseteq \Lambda \cap S^{d-1}$ of size C^d such that if M preserves the norm of every $x \in T$ then M preserves all of Λ up to $1 + \epsilon$. See [28] for details.

► **Lemma 14** (Khinchine's inequality [10]). *Let $p \geq 1$, $x \in \mathbb{R}^d$, and $\sigma_i \mathbb{R}^d$ be independent Rademacher ± 1 random variables. Then*

$$\left\| \sum_{i=1}^d \sigma_i x_i \right\|_p \lesssim \sqrt{p} \|x\|_2.$$

3 Technical Overview

The main component of any tensor sketch is a matrix $M : \mathbb{R}^{m \times d_1 d_2}$ such that $\|Mx\|_2 \approx \|x\|_2$ and which can be applied efficiently (faster than $md_1 d_2$) to simple tensors $x = x^{(1)} \otimes x^{(2)}$, where $x^{(1)} \in \mathbb{R}^{d_1}, x^{(2)} \in \mathbb{R}^{d_2}$.

If $x^{(1)}$ and $x^{(2)}$ are ± 1 vectors, sampling from x works well and can be done without actually constructing x . For this reason a natural general sketch is $S(M^{(1)}x^{(1)} \otimes M^{(2)}x^{(2)})$, where $M^{(1)}$ and $M^{(2)}$ are random rotations.

The original Tensor Sketch did $\mathcal{F}^{-1}(\mathcal{F}C^{(1)}x \circ \mathcal{F}C^{(2)}x)$, where $C^{(1)}$ and $C^{(2)}$ are Count Sketch matrices. At first sight this may look somewhat different, but we can ignore the orthonormal \mathcal{F}^{-1} , and then we have $M^{(1)}x^{(1)} \circ M^{(2)}x^{(2)}$ which is just sampling the diagonal of $M^{(1)}x^{(1)} \otimes M^{(2)}x^{(2)}$. Since $M^{(1)}$ and $M^{(2)}$ are independent, sampling the diagonal works as well as any other subset of the same size.

Since Tensor Sketch only used 2nd moment analysis, the natural technical question is “how well does $M^{(1)}x^{(1)} \circ M^{(2)}x^{(2)}$ really work?” when $M^{(1)}$ and $M^{(2)}$ can be anything. In Theorem 27 we show that an embedding dimension of $m = \Theta(\epsilon^{-2} \log 1/\delta + \epsilon^{-1} (\log 1/\delta)^2)$ is both sufficient and necessary for (sub)-gaussian matrices, which we conjecture is optimal across all distributions.

Sub-gaussian matrices however still take md time to evaluate, so our tensor sketch would still take $m(d_1 + d_2)$ time in total. We really want $M^{(1)}$ and $M^{(2)}$ to have fast matrix-vector multiplication. It is thus natural to analyze the above scheme where $M^{(1)}$ and $M^{(2)}$ are Fast Johnson Lindenstrauss matrices ala [1, 15]. We do this in Section 5 and show that $m = \epsilon^{-2} (\log 1/\delta) (\log 1/\epsilon\delta)^2$ suffices. For ϵ not too small, this matches our suggested optimum by one $\log 1/\delta$ factor.

The final challenge is to scale up to larger tensors than order 2. Our Lemma 19 shows and exponential dependency: $m = \epsilon^{-2} (\log 1/\delta) (\log 1/\epsilon\delta)^c$, which would be rather unfortunate. Luckily it turns out, that by continuously ‘squashing’ the dimension back down to $\epsilon^{-2} (\log 1/\delta) (\log 1/\epsilon\delta)^2$, we can avoid this explosion.

While we usually think of applying our sketching matrix to simple tensors, we always analyze everything assuming the input has full rank. This adds some extra difficulty to the analysis, but it is worth it, since by showing that our matrix has the so called JL-moment-property, we get that it is also a subspace embedding for free, by Lemma 12 and Lemma 13.

4 The High Probability Tensor Sketch

In this section we will prove Theorem 3 which is the backbone of our theorems. Theorem 2 will follow as an easy corollary, while Theorem 1 is completed in the next section.

Before we show the full theorem we will consider a slightly easier construction. Given independent random matrices $Q^{(2)}, \dots, Q^{(c)} \in \mathbb{R}^{m \times dm}$, from a distribution to be discussed later, and $Q^{(1)} \in \mathbb{R}^{m \times d}$, we define $M^{(0)} = 1 \in \mathbb{R}$ and recursively $M^{(k)} = Q^{(k)}(M^{(k-1)} \otimes I_d)$

XX:8 High Probability Tensor Sketch

237 for $k \in [c]$. The goal of this section is to show that M^k has JL- and related properties when the
 238 $Q^{(i)}$ s have, and that $M^{(k)}$ can be evaluated efficiently on simple tensors, $x^{(1)} \otimes \dots \otimes x^{(k)} \in \mathbb{R}^{d^k}$,
 239 for $k \in [c]$.

240 First we show a rather simple fact which will prove to be quite powerful.

241 ► **Lemma 15.** *Let $\varepsilon \in (0, 1)$ and $\delta > 0$. If $P \in \mathbb{R}^{m_1 \times d_1}$ and $Q \in \mathbb{R}^{m_2 \times d_2}$ are two matrices
 242 with (ε, δ) -JL moment property, then $P \oplus Q \in \mathbb{R}^{(m_1+m_2) \times (d_1+d_2)}$ has (ε, δ) -JL moment
 243 property.*

244 **Proof.** Let $x \in \mathbb{R}^{d_1+d_2}$ and choose $y \in \mathbb{R}^{d_1}$ and $z \in \mathbb{R}^{d_2}$ such that $x = y \oplus z$. Now using the
 245 triangle inequality and JL moment property, we get that

$$\begin{aligned} 246 \quad \left| \|(P \oplus Q)x\|_2^2 - \|x\|_2^2 \right|_p &\leq \left| \|Py\|_2^2 - \|y\|_2^2 \right|_p + \left| \|Qz\|_2^2 - \|z\|_2^2 \right|_p \\ 247 \quad &\leq \varepsilon \delta^{1/p} \|y\|_2^2 + \varepsilon \delta^{1/p} \|z\|_2^2 \\ 248 \quad &= \varepsilon \delta^{1/p} \|x\|_2^2, \end{aligned}$$

250 since $\|y\|_2^2 + \|z\|_2^2 = \|y \oplus z\|_2^2$ by disjointness. ◀

251 An easy consequence of this lemma is that for any matrix T , $I_\ell \otimes T$ has (ε, δ) -JL moment
 252 property when T has (ε, δ) -JL moment property, since $I_\ell \otimes Q = \underbrace{Q \oplus Q \oplus \dots \oplus Q}_{\ell \text{ times}}$.

253 Similarly, $Q \otimes I_\ell$ has (ε, δ) -JL moment property, since you can obtain $Q \otimes I_\ell$ by reordering
 254 the rows of $I_\ell \otimes Q$, which trivially does not change the JL moment property.

255 It is now easy to show that $M^{(k)}$ has JL-property when $Q^{(1)}, \dots, Q^{(k)}$ has JL-property
 256 for $k \in [c]$.

257 ► **Lemma 16.** *Let $\varepsilon \in (0, 1)$ and $\delta > 0$. If $Q^{(1)}, \dots, Q^{(c)}$ has the $(\varepsilon/2c, \delta/c)$ -JL property,
 258 then $M^{(k)}$ has $(k/c\varepsilon, k/c\delta)$ -jl property for every $k \in [c]$.*

259 **Proof.** Let $k \in [c]$ be fixed. We note that an alternative way of expressing $M^{(k)}$ is as follows:

$$261 \quad M^{(k)} = Q^{(k)}(Q^{(k-1)} \otimes I_d)(Q^{(k-2)} \otimes I_{d^2}) \dots (Q^{(1)} \otimes I_{d^{k-1}})$$

262 Let $x \in \mathbb{R}^{d^k}$ be any vector. Define $x^{(i)} = (Q^{(i)} \otimes I_{d^{k-i}})x^{(i-1)}$ for $i \in [k]$ and $x^{(0)} = x$.
 263 Since $Q^{(i)}$ has $(\varepsilon/2k, \delta/k)$ -JL property then $Q^{(i)} \otimes I_{d^i}$ has $(\varepsilon/2c, \delta/c)$ -JL property by the
 264 previous discussion, hence $\Pr \left[\left| \|x^{(i)}\|_2^2 - \|x^{(i)}\|_2^2 \right| \geq \varepsilon/2c \|x^{(i)}\|_2^2 \right] \leq \delta/c$. Now a simple union
 265 bound give us that

$$266 \quad 1 - k/c\varepsilon \leq (1 - \varepsilon/2c)^k \leq \left| \|Mx\|_2^2 - \|x\|_2^2 \right| \leq (1 + \varepsilon/2c)^k \leq 1 + k/c\varepsilon$$

267 with probability at least $1 - k/c\delta$, which finishes the proof. ◀

268 ► **Corollary 17.** *Let $\varepsilon \in (0, 1)$. If $Q^{(1)}, \dots, Q^{(c)}$ has the $(\varepsilon/(2\lambda c), \delta/c)$ -JL property, then
 269 $M^{(k)}$ is a λ -subspace embedding.*

270 **Proof.** This follows from lemma 12. ◀

271 Note that if $Q^{(i)}x$, where $x \in \mathbb{R}^{d^k}$, can be evaluated in time T for every $i \in [c] \setminus \{1\}$, and
 272 $Q^{(1)}y$, where $y \in \mathbb{R}^d$, also can be evaluated in time T' , then $M^{(k)}z$, where $z \in \mathbb{R}^{d^k}$, can be
 273 evaluated in time $T'd^{k-1} + T \sum_{i=0}^{k-2} d^i = T'd^{k-1} + T(d^{k-1} - 1)/(d - 1) = \Theta(T'd^{k-1} + Td^{k-2})$.
 274 Meanwhile, if $x \in \mathbb{R}^{d^k}$ is on the form $x^{(1)} \otimes \dots \otimes x^{(k)}$, we have $M^{(k)}x = Q^{(k)}(M^{(k-1)}(x^{(1)} \otimes$

275 $\dots x^{(k-1)}) \otimes x^{(k)}$. Now an easy induction argument shows that this allows evaluation in
 276 time $O(T' + Tk)$, which is exponentially faster.

277 Using this construction it now becomes easy to sketch polynomials. More precisely, let
 278 $t \in \mathbb{Z}_{>0}$, $k_i \in [c]$ for every $i \in [t]$, then the matrix $M = \bigoplus_{i \in [t]} M^{(k_i)}$ has (ε, δ) -JL property
 279 and can be evaluated at the vector $x = \bigoplus_{i \in [t]} \bigotimes_{j \in [k_i]} x^{(i,j)}$ in time $O(T't + T \sum_{i \in [t]} k_i)$,
 280 where $x^{(i,j)} \in \mathbb{R}^d$ for every $i \in [t], j \in [k_i]$.

281 This discussion proves Theorem 3. Note that if we apply a Fast Johnson Lindenstrauss
 282 Transform between every direct sum we can obtain an output dimension of $O(m)$.

283 **► Example 18.** Often it is possible to get an even faster evaluation time if the input has
 284 even more structure. For example consider the matrix $M = \bigoplus_{i \in [c]} M^{(i)}$ and the vector
 285 $z = \bigoplus_{i \in [c]} \bigotimes_{j \in [i]} x^{(j)}$, where $x^{(j)} \in \mathbb{R}^d$ for every $j \in [c]$. Then Mz can be evaluated in time
 286 $O(T' + cT)$ by exploiting the fact that

$$287 \quad M^{(k)} \left(\bigotimes_{j \in [k]} x^{(j)} \right) = Q^{(k)} \left(M^{(k-1)} \left(\bigotimes_{j \in [k-1]} x^{(j)} \right) \otimes x^{(k)} \right),$$

288 so we can use the previous calculations.

289 As promised we now get the proof of Theorem 2 by choosing $Q^{(1)}, \dots, Q^{(c)}$ to be Fast
 290 Johnson Lindenstrauss Matrices. Using the analysis from Kraher et al. [15] they can be
 291 evaluated in time $O(md \log md)$ and if we set $m = \tilde{O}(c^2 \log(1/\delta)/\varepsilon^2)$ then $Q^{(1)}, \dots, Q^{(c)}$ has
 292 $(\varepsilon/2c, \delta/c)$ -JL property. Now Theorem 3 give us the result.

293 5 Fast Constructions

294 The purpose of this section is to show the following lemma:

295 **► Lemma 19.** Let $c \in \mathbb{Z}_{>0}$, and $(D^{(i)})_{i \in [c]} \in \prod_{i \in [c]} \mathbb{R}^{d_i \times d_i}$ be independent diagonal matrices
 296 with independent Rademacher variables. Define $d = \prod_{i \in [c]} d_i$ and $D = \bigotimes_{i \in [c]} D_i \in \mathbb{R}^{d \times d}$.
 297 Let $S \in \mathbb{R}^{m \times d}$ be an independent sampling matrix which samples exactly one coordinate per
 298 row. Let $x \in \mathbb{R}^d$ be any vector and $p \geq 1$, then

$$299 \quad \left\| \frac{1}{m} \|SHDx\|_2^2 - \|x\|_2^2 \right\|_p \lesssim \sqrt{p} (p + \log m)^{c/2} \|x\|_2^2 / \sqrt{m} + p (p + \log m)^c \|x\|_2^2 / m.$$

301 Setting $m = O(\varepsilon^{-2} \log 1/\delta (\log 1/\varepsilon \delta)^c)$ thus suffices for SHD to have the (ε, δ) -JL-moment-
 302 property. This then gives (by lemma 9 and 12) that SHD is a subspace embedding.

303 We note that SHD can be applied efficiently to simple tensors by the relation:

$$304 \quad SH_{d_1 d_2} (D^{(1)} \otimes D^{(2)}) (x \otimes y) = (S^{(1)} \otimes S^{(2)}) (H_{d_1} \otimes H_{d_2}) (D^{(1)} \otimes D^{(2)}) (x \otimes y) \\ 305 \quad = S^{(1)} H_{d_1} D^{(1)} x \circ S^{(2)} H_{d_2} D^{(2)} y,$$

307 where H_n is the size n Hadamard matrix and $S^{(1)}$ and $S^{(2)}$ are independent sampling matrices.
 308 Combining this fact with the construction in the previous section gives Theorem 1.

309 The rest of this section is devoted to proving Lemma 19. We first show two technical
 310 lemmas, which seem like they could be useful for many other things.

311 **► Lemma 20.** Let $p \geq 1$, $c \in \mathbb{Z}_{>0}$, and $(\sigma^{(i)})_{i \in [c]} \in \prod_{i \in [c]} \mathbb{R}^{d_i}$ be independent Rademacher
 312 vectors. Let $a_{i_0, \dots, i_{c-1}} \in \mathbb{R}$ for every $i_j \in [d_j]$ and every $j \in [c]$, then

$$313 \quad \left\| \sum_{i_1 \in [d_1], \dots, i_c \in [d_c]} \prod_{j \in [c]} \sigma_{i_j}^{(j)} a_{i_0, \dots, i_{c-1}} \right\|_p \lesssim p^{c/2} \left(\sum_{i_1 \in [d_1], \dots, i_c \in [d_c]} a_{i_0, \dots, i_{c-1}}^2 \right)^{1/2} = p^{c/2} \|a\|_{HS}.$$

314

XX:10 High Probability Tensor Sketch

315 **Proof.** The proof will be by induction on c . For $c = 1$ then the result is just Khintchine's
 316 inequality (Lemma 14). So assume that the result is true for every value up to c . Using the
 317 induction hypothesis we get that

$$\begin{aligned}
 318 \quad & \left\| \sum_{\substack{i_1 \in [d_1], j \in [c] \\ \dots, i_c \in [d_c]}} \prod \sigma_{i_j}^{(j)} a_{i_1, \dots, i_c} \right\|_p = \left\| \sum_{\substack{i_1 \in [d_1], \\ \dots, i_{c-1} \in [d_{c-1}]}} \prod_{j \in [c-1]} \sigma_{i_j}^{(j)} \left(\sum_{i_c \in [d_c]} \sigma_{i_c}^{(c)} a_{i_1, \dots, i_c} \right) \right\|_p \\
 319 \quad & \lesssim p^{(c-1)/2} \left\| \left(\sum_{\substack{i_1 \in [d_1], \\ \dots, i_{c-1} \in [d_{c-1}]} \left(\sum_{i_c \in [d_c]} \sigma_{i_c}^{(c)} a_{i_1, \dots, i_c} \right)^2 \right)^{1/2} \right\|_p \quad (\text{I.H.}) \\
 320 \quad & = p^{(c-1)/2} \left\| \sum_{\substack{i_1 \in [d_1], \\ \dots, i_{c-1} \in [d_{c-1}]} \left(\sum_{i_c \in [d_c]} \sigma_{i_c}^{(c)} a_{i_1, \dots, i_c} \right)^2 \right\|_{p/2}^{1/2} \\
 321 \quad & \leq p^{(c-1)/2} \left(\sum_{\substack{i_1 \in [d_1], \\ \dots, i_{c-1} \in [d_{c-1}]} \left\| \left(\sum_{i_c \in [d_c]} \sigma_{i_c}^{(c)} a_{i_1, \dots, i_c} \right)^2 \right\|_{p/2} \right)^{1/2} \quad (\text{Triangle}) \\
 322 \quad & = p^{(c-1)/2} \left(\sum_{\substack{i_1 \in [d_1], \\ \dots, i_{c-1} \in [d_{c-1}]} \left\| \sum_{i_c \in [d_c]} \sigma_{i_c}^{(c)} a_{i_1, \dots, i_c} \right\|_p^2 \right)^{1/2} \\
 323 \quad & \lesssim p^{c/2} \left(\sum_{i_1 \in [d_1], \dots, i_c \in [d_c]} a_{i_1, \dots, i_c}^2 \right)^{1/2} \quad (\text{Khintchine}) \\
 324
 \end{aligned}$$

325 where the last inequality is by using Khintchine's inequality. Plugging this into the previous
 326 inequality finishes the induction step and hence the proof. ◀

327 The next lemma we need is a type of Chernoff bound for p th moments.

328 ▶ **Lemma 21.** *Let $p \geq 2$ and X_0, \dots, X_{k-1} be independent non-negative random variables
 329 with p -moment, then*

$$330 \quad \left\| \sum_{i \in [k]} (X_i - \mathbb{E}[X_i]) \right\|_p \lesssim \sqrt{p} \sqrt{\sum_{i \in [k]} \mathbb{E}[X_i]} \left\| \max_{i \in [k]} X_i \right\|_p^{1/2} + p \left\| \max_{i \in [k]} X_i \right\|_p$$

Proof.

$$\begin{aligned}
 331 \quad & \left\| \sum_{i \in [k]} (X_i - \mathbb{E}[X_i]) \right\|_p \lesssim \left\| \sum_{i \in [k]} \sigma_i X_i \right\|_p \quad (\text{Symmetrization}) \\
 332 \quad & \lesssim \sqrt{p} \left\| \sqrt{\sum_{i \in [k]} X_i^2} \right\|_p \quad (\text{Khintchine's inequality}) \\
 333 \quad & = \sqrt{p} \left\| \sum_{i \in [k]} X_i^2 \right\|_{p/2}^{1/2} \\
 334 \quad & \leq \sqrt{p} \left\| \max_{i \in [k]} X_i \right\|_p^{1/2} \left\| \sum_{i \in [k]} X_i \right\|_p^{1/2} \quad (\text{H\"older's inequality}) \\
 335 \quad & \leq \sqrt{p} \left\| \max_{i \in [k]} X_i \right\|_p^{1/2} \sqrt{\sum_{i \in [k]} \mathbb{E}[X_i]} \\
 336 \quad & \quad + \sqrt{p} \left\| \max_{i \in [k]} X_i \right\|_p^{1/2} \left\| \sum_{i \in [k]} (X_i - \mathbb{E}[X_i]) \right\|_p^{1/2} \quad (\text{Triangle inequality}) \\
 337
 \end{aligned}$$

338 Now let $C = \left\| \sum_{i \in [k]} (X_i - \mathbb{E}[X_i]) \right\|_p^{1/2}$, $B = \sqrt{\sum_{i \in [k]} \mathbb{E}[X_i]}$, and $A = \sqrt{p} \left\| \max_{i \in [k]} X_i \right\|_p^{1/2}$. then
 339 we have shown $C^2 \leq A(B + C)$. That implies C is smaller than the largest of the roots of
 340 the quadratic. Solving this quadratic inequality gives $C^2 \lesssim AB + A^2$ which is the result. ◀

341 We can finally go ahead and prove Lemma 19.

342 **Proof.** For every $i \in [m]$ we let S_i be the random variable that says which coordinate the i 'th
343 row of S samples, and we define the random variable $Z_i = M_i x_i = H_{S_i} D x_i$. We note that
344 since the variables $(S_i)_{i \in [m]}$ are independent then the variables $(Z_i)_{i \in [m]}$ are conditionally
345 independent given D , that is, if we fix D then $(Z_i)_{i \in [m]}$ are independent.

346 Using Lemma 21 we get that

$$347 \quad \left\| \frac{1}{m} \sum_{i \in [m]} Z_i^2 - \|x\|_2^2 \right\|_p \lesssim \sqrt{p} \left(\frac{1}{m} \sum_{i \in [m]} \mathbb{E}[Z_i^2 | D] \right)^{1/2} \left\| \max_{i \in [m]} \frac{1}{m} Z_i^2 \right\|_p^{1/2} + p \left\| \max_{i \in [m]} \frac{1}{m} Z_i^2 \right\|_p$$

348 (1)

349 It follows easily that $\mathbb{E}[Z_i^2 | D] = \|x\|_2^2$ from the fact that $\|HDx\|_2^2 = d\|x\|_2^2$, hence
350 $\left(\frac{1}{m} \sum_{i \in [m]} \mathbb{E}[Z_i^2 | D] \right)^{1/2} = \|x\|_2$. Now we just need to bound $\left\| \max_{i \in [m]} \frac{1}{m} Z_i^2 \right\|_p =$
351 $\frac{1}{m} \left\| \max_{i \in [m]} Z_i^2 \right\|_p$. First we note that

$$352 \quad \left\| Z_i^2 \right\|_p = \left\| (H_{S_i} D x_i)^2 \right\|_p \lesssim p^c \|x\|_2^2$$

353

354 by Khintchine's inequality. Let $q = \max\{p, \log m\}$, then we get that

$$355 \quad \left\| \max_{i \in [m]} Z_i^2 \right\|_p \leq \left\| \max_{i \in [m]} Z_i^2 \right\|_q \leq \left(\sum_{i \in [m]} \|Z_i^2\|_q \right)^{1/q} \leq m^{1/q} q^c \|x\|_2^2$$

356

357 Now since $q \geq \log m$ then $m^{1/q} \leq 2$ so $\left\| \max_{i \in [m]} Z_i^2 \right\|_p \lesssim q^c \|x\|_2^2 \leq (p + \log m)^c \|x\|_2^2$.
358 Plugging this into 1 finishes the proof. ◀

359 6 Applications

360 The classic application of Tensor Sketching is compact bilinear pooling (or multilinear). This
361 simply corresponds to expanding x and $x^{\otimes 2}$ (bilinear pooling) and then hashing back to a smaller
362 size (compact). First discussed in [9] which showed how to do back-propagation through
363 a tensor-sketch layer. Then applied to all many applications such as visual convolutional
364 models [20], question answering [8], visual reasoning [12], video classification [21].

365 These results are usually given without any particular guarantees, but we can also use
366 polynomial embeddings for concrete algorithms using the following lemma:

367 6.1 Sketching Polynomials

368 ▶ **Theorem 22.** *Given any degree c polynomial, $P(z) = \sum_{i=0}^c a_i z^i$, there are a pair of*
369 *embeddings $f, g: \mathbb{R}^d \rightarrow \mathbb{R}^m$, such that for any $x, y \in \mathbb{R}^d$, the inner product*

$$370 \quad \langle f(x), g(y) \rangle = (1 \pm \epsilon) P(\langle x, y \rangle)$$

371

372 *with probability at least $1 - \delta$. Using Construction A we set $m = O(c^2 \epsilon^{-2} (\log 1/\delta) (\log 1/\epsilon \delta)^2)$,*
373 *and $f(x)$ and $g(y)$ can be computed in $O(c(d \log d + m \log m))$ time. Using Construction B*
374 *we set $m = \tilde{O}(c^2 \epsilon^{-2} (\log 1/\delta))$, and $f(x)$ and $g(y)$ can be computed in $O(cm \min(d, m))$ time.*

375 **Proof.** We note that

$$376 \quad P(\langle x, y \rangle) = \sum_{i=0}^c a_i \langle x, y \rangle^i = \sum_{i=0}^c \langle a_i x^{\otimes i}, y^{\otimes i} \rangle = \left\langle \bigoplus_{i=0}^c a_i x^{\otimes i}, \bigoplus_{i=0}^c y^{\otimes i} \right\rangle.$$

377

378 So using Theorem 3 together with Construction A and Construction B give the result. ◀

XX:12 High Probability Tensor Sketch

379 We note that the output dimension m in the theorem can be improved by applying a Fast
380 Johnson Lindenstrauss Transform in the end.

381 ► **Example 23 (Explicit Gaussian Kernel).** Say we want $\langle f(x), g(y) \rangle \approx \exp(-\langle x, y \rangle^2) =$
382 $\sum_{k=0}^c (-1)^k \langle x, y \rangle^{2k} / k! + O(\langle x, y \rangle^{2c+1} / c!).$ Then using Theorem 22 we can obtain $\langle f(x), g(y) \rangle =$
383 $(1 \pm \varepsilon) \sum_{k=0}^c (-1)^k \langle x, y \rangle^{2k} / k!,$ hence get an approximation of $\varepsilon + O(\langle x, y \rangle^{2c+1} / c!)$ with prob-
384 ability $1 - \delta,$ using $O(c(d \log d + m \log m))$ time where $m = O(c^3 \varepsilon^{-2} (\log 1/\delta) (\log 1/\varepsilon \delta)^2).$

385 6.2 Embeddings

386 ► **Lemma 24 (Symmetric Polynomials).** *Given any degree polynomial, $P(x_1, \dots, x_d, y_1, \dots, y_d)$*
387 *with k monomials we can make an embedding $f, g : R^d \rightarrow R^m$ such that*

$$388 \quad E \langle f(x), g(y) \rangle = \sum_{\pi} P(x_{\pi(1)}, \dots, x_{\pi(d)}, y_{\pi(1)}, \dots, y_{\pi(d)}).$$

390 $f(x)$ and $g(y)$ can be computed in time $k 4^{\Delta}$ sketching operations, where Δ is the maximum
391 combined degree of a monomial. (E.g. 4 for $x_1^2 y_1 y_2$.)

392 For x and y boolean, we get $\|f(x)\|_2^2 \leq k(\|x\|_2^2 + \kappa - 1)! / (\|x\|_2^2 - 1)! \leq k(\|x\|_2^2 + \kappa - 1)^{\kappa}.$
393 where κ is the number of different indicies in a monomial. E.g. for $x_1 y_1 x_2$ it is 2.

394 **Proof.** See the Appendix Section 7.2. ◀

395 ► **Example 25 (Triangle counting).** Say you have a database of graphs, $\mathcal{G},$ seen as binary
396 vectors in $\{0, 1\}^{\varepsilon}.$ Given a query graph, $G,$ you want to find $G' \in \mathcal{G}$ such that the number of
397 triangles in $G \cup G'$ is maximized.

398 We construct the polynomial $P(x, y) = (x_1 + y_1 - x_1 y_1)(x_2 + y_2 - x_2 y_2)(x_3 + y_3 - x_3 y_3)$
399 which is 1 exactly when $x \cup u$ has a triangle on edges 1, 2, 3. The maximum number of
400 different indicies is $\kappa = 3.$ The maximum number of triangles (with ordering) is $6 \binom{d, 3}{} d^3.$
401 We have $\|f(x)\|_2 \|g(x)\| \approx (d + 3)^3,$ so to get precision within 1% of the maximum number,
402 we need to set $\epsilon < 0.01 d^3 / (d + 3)^3.$

403 Since our approximation works with high probability, we can take a union bound and
404 plug the embedded vectors into the standard data structure of [3] or others.

405 6.3 Oblivious Subspace Embeddings

406 In [28] the authors show a number of applications of polynomial kernels in oblivious subspace
407 embeddings. They also show that the original tensor sketch [25] is an oblivious subspace
408 embedding when the sketch matches the size described in the introduction. It is shown how
409 each of:

- 410 1. Approximate Kernel PCA and Low Rank Approximation,
- 411 2. Regularizing Learning With the Polynomial Kernel,
- 412 3. Approximate Kernel Principal Component Regression,
- 413 4. Approximate Kernel Canonical Correlation Analysis,

414 can be computed with state of the art performance.

415 However, each of the applications encounter an exponential dependency on $c.$ They also
416 inherit the tensor-sketch linear dependency on the inverse error probability. Our sketch
417 improves each of these aspects directly black box.

418 — References —

- 419 1 Nir Ailon and Bernard Chazelle. Approximate nearest neighbors and the fast johnson-
420 lindenstrauss transform. In *Proceedings of the thirty-eighth annual ACM symposium on Theory*
421 *of computing*, pages 557–563. ACM, 2006.
- 422 2 Josh Alman, Timothy M Chan, and Ryan Williams. Polynomial representations of threshold
423 functions and algorithmic applications. In *2016 IEEE 57th Annual Symposium on Foundations*
424 *of Computer Science (FOCS)*, pages 467–476. IEEE, 2016.
- 425 3 Alexandr Andoni, Thijs Laarhoven, Ilya Razenshteyn, and Erik Waingarten. Optimal hashing-
426 based time-space trade-offs for approximate near neighbors. In *Proceedings of the Twenty-Eighth*
427 *Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 47–66. Society for Industrial
428 and Applied Mathematics, 2017.
- 429 4 Anonymous. Sketching high-degree polynomial kernels. In double blind review, 2019.
- 430 5 Haim Avron, Huy L. Nguyen, and David P. Woodruff. Subspace embeddings for the
431 polynomial kernel. In *Advances in Neural Information Processing Systems 27: An-*
432 *ual Conference on Neural Information Processing Systems 2014, December 8-13 2014,*
433 *Montreal, Quebec, Canada*, pages 2258–2266, 2014. URL: [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/5240-subspace-embeddings-for-the-polynomial-kernel)
434 [5240-subspace-embeddings-for-the-polynomial-kernel](http://papers.nips.cc/paper/5240-subspace-embeddings-for-the-polynomial-kernel).
- 435 6 Moses S Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings*
436 *of the thirty-fourth annual ACM symposium on Theory of computing*, pages 380–388. ACM,
437 2002.
- 438 7 Michael B Cohen, TS Jayram, and Jelani Nelson. Simple analyses of the sparse johnson-
439 lindenstrauss transform. In *1st Symposium on Simplicity in Algorithms (SOSA 2018)*. Schloss
440 Dagstuhl-Leibniz-Zentrum fuer Informatik, 2018.
- 441 8 Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus
442 Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual
443 grounding. *arXiv preprint arXiv:1606.01847*, 2016.
- 444 9 Yang Gao, Oscar Beijbom, Ning Zhang, and Trevor Darrell. Compact bilinear pooling. In
445 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 317–326,
446 2016.
- 447 10 Uffe Haagerup and Magdalena Musat. On the best constants in noncommutative khintchine-
448 type inequalities. *Journal of Functional Analysis*, 250(2):588–624, 2007.
- 449 11 Piotr Indyk and Assaf Naor. Nearest-neighbor-preserving embeddings. *ACM Transactions on*
450 *Algorithms (TALG)*, 3(3):31, 2007.
- 451 12 Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick,
452 and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary
453 visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
454 *Recognition*, pages 2901–2910, 2017.
- 455 13 William B Johnson and Joram Lindenstrauss. Extensions of lipschitz mappings into a hilbert
456 space. *Contemporary mathematics*, 26(189-206):1, 1984.
- 457 14 William B Johnson, Gideon Schechtman, and Joel Zinn. Best constants in moment inequalities
458 for linear combinations of independent and exchangeable random variables. *The Annals of*
459 *Probability*, pages 234–253, 1985.
- 460 15 Felix Kraher and Rachel Ward. New and improved johnson–lindenstrauss embeddings via
461 the restricted isometry property. *SIAM Journal on Mathematical Analysis*, 43(3):1269–1281,
462 2011.
- 463 16 Kasper Green Larsen and Jelani Nelson. Optimality of the johnson-lindenstrauss lemma. In
464 *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages
465 633–638. IEEE, 2017.
- 466 17 Rafał Łatała et al. Estimates of moments and tails of gaussian chaoses. *The Annals of*
467 *Probability*, 34(6):2315–2331, 2006.

- 468 **18** Quoc Viet Le, Tamás Sarlós, and Alexander Johannes Smola. Fastfood: Approximate kernel
 469 expansions in loglinear time. *CoRR*, abs/1408.3060, 2014. URL: [http://arxiv.org/abs/](http://arxiv.org/abs/1408.3060)
 470 [1408.3060](http://arxiv.org/abs/1408.3060), [arXiv:1408.3060](https://arxiv.org/abs/1408.3060).
- 471 **19** Joseph Lehec. Moments of the gaussian chaos. In *Séminaire de Probabilités XLIII*, pages
 472 327–340. Springer, 2011.
- 473 **20** Tsung-Yu Lin, Aruni RoyChowdhury, and Subhransu Maji. Bilinear cnn models for fine-grained
 474 visual recognition. In *Proceedings of the IEEE international conference on computer vision*,
 475 pages 1449–1457, 2015.
- 476 **21** Antoine Miech, Ivan Laptev, and Josef Sivic. Learnable pooling with context gating for video
 477 classification. *arXiv preprint arXiv:1706.06905*, 2017.
- 478 **22** Cameron Musco and Christopher Musco. Recursive sampling for the nystrom method. In
 479 *Advances in Neural Information Processing Systems*, pages 3833–3845, 2017.
- 480 **23** Nam H Nguyen, Petros Drineas, and Trac D Tran. Tensor sparsification via a bound on
 481 the spectral norm of random tensors. *Information and Inference: A Journal of the IMA*,
 482 4(3):195–229, 2015.
- 483 **24** Rasmus Pagh. Compressed matrix multiplication. *ACM Transactions on Computation Theory*
 484 (*TOCT*), 5(3):9, 2013.
- 485 **25** Ninh Pham and Rasmus Pagh. Fast and scalable polynomial kernels via explicit feature maps.
 486 In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery*
 487 *and data mining*, pages 239–247. ACM, 2013.
- 488 **26** Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances*
 489 *in neural information processing systems*, pages 1177–1184, 2008.
- 490 **27** Gregory Valiant. Finding correlations in subquadratic time, with applications to learning
 491 parities and the closest pair problem. *Journal of the ACM (JACM)*, 62(2):13, 2015.
- 492 **28** David P. Woodruff. Sketching as a tool for numerical linear algebra. *Foundations and Trends*
 493 *in Theoretical Computer Science*, 10(1-2):1–157, 2014. URL: [https://doi.org/10.1561/](https://doi.org/10.1561/04000000060)
 494 [04000000060](https://doi.org/10.1561/04000000060), [doi:10.1561/04000000060](https://doi.org/10.1561/04000000060).

495 **7** Appendix

496 **7.1** Subgaussian construction

497 Before stating the theorem, we not the following matrix product:

498 **► Definition 26** (Face-splitting product). *Defined between to matrices as the Kronecker-product*
 499 *between pairs of rows:*

$$500 \quad C \bullet D = \begin{bmatrix} C_1 \otimes D_1 \\ C_2 \otimes D_2 \\ \dots \\ C_n \otimes D_n \end{bmatrix}.$$

502 Face-splitting product has the relation $(A \bullet B)(x \otimes y) = Ax \circ By$.

503 **► Theorem 27** (Subgaussian). *Let $T, M \in \mathbb{R}^{m \times d}$ have iid. sub-gaussian entries, then*
 504 *$\| \frac{1}{\sqrt{m}}(T \bullet M)x \|_2^2 - \|x\|_2^2 \|_p \leq \sqrt{p/m} + pq/m$, where $q = \max(p, m)$.*

505 This immediately implies that for $m = \Omega(\epsilon^{-2} \log 1/\delta + \epsilon^{-1}(\log 1/\delta)(\log 1/\epsilon\delta))$, $T \bullet M$ has
 506 the JL-moment property.

507 We note that the analysis is basically optimal. Assume M and T were iid. Gaussian
 508 matrices and $x = e_1^{\otimes 2}$ were a simple tensor with a single 1 entry. Then $\| \frac{1}{\sqrt{m}}(M \bullet T)x \|_2^2 - \|x\|_2^2 \|_p =$
 509 $\| \|Mx' \circ Tx'\|_2^2 - 1 \| \sim |(gg')^2 - 1|$ for $g, g' \in R$ iid. Gaussians. Now $\Pr[(gg')^2 > (1 + \epsilon)] \approx$
 510 $\exp(-\min(\epsilon, \sqrt{\epsilon}))$, thus requiring $m = \Omega(\epsilon^{-2} \log 1/\delta + \epsilon^{-1}(\log 1/\delta)^2)$ matching our bound
 511 up to a $\log 1/\epsilon$.

512 **Proof.** Let $Q = T \bullet M \in \mathbb{R}^{n \times ab}$.

$$\begin{aligned} 513 \quad \|(T \bullet M)x\|_2^2 &= \sum_k ((T \bullet M)_{k,x})^2 \\ 514 \quad &= \sum_k (TU_{(i)})_{k,k}^2 \\ 515 \end{aligned}$$

516 where $U_{(i)} = (XM^T)_{(i)} = (M_i X)^T$, when x is seen as a $d \times d$ matrix.

517 We then have sub-gaussians:

$$\begin{aligned} 518 \quad E\|U_{(i)}\|_2^2 &= E\|M_i X\|_2^2 = \sum_k E(M_i X_{(k)})^2 = \sum_k \|X_{(k)}\|^2 = \|X\|_F^2 = 1 \\ 519 \quad \|\|U_{(i)}\|_2\|_p &\leq \|\|M_i X^T\|_2\|_p + 1 \leq \sqrt{p}\|X^T X\|_F + p\|X^T X\| + 1 \leq p \\ 520 \quad \|\sum_k \|U_{(k)}\|_2^4\|_p &\leq \sqrt{pk} + k + pq^2, \\ 521 \end{aligned}$$

522 The last bound followed from independence and bounded variance of the $\|U_{(k)}\|_2$ s. It is
523 possible to go without this assumption though, suffering a small loss in the final dimension.

524 We bound

$$525 \quad \|(T \bullet M)x\|_2^2/k - \|x\|_2^2\|_p \leq \|\sum_k (T_k U_{(k)})^2 - \|U\|_F^2\|_p/k \quad (2)$$

$$526 \quad + \|\|U\|_F^2/k - \|x\|_2^2\|_p. \quad (3)$$

528 Bounding (3) follows simply from the JL property of M . Bounding (2) is a bit trickier, and
529 we use the Hanson-Wright inequality:

$$\begin{aligned} 530 \quad \|\sum_k (T_k U_{(k)})^2 - \|U\|_F^2\|_p &\leq \|\sqrt{p}(\sum_k \|U_{(k)} U_{(k)}^T\|_2^2)^{1/2} + p \max_k \|U_{(k)} U_{(k)}^T\|_p \\ 531 \quad &\leq \|\sqrt{p}(\sum_k \|U_{(k)}\|_2^4)^{1/2} + p \max_k \|U_{(k)}\|_2^2\|_p \\ 532 \quad &\leq \sqrt{p} \|\sum_k \|U_{(k)}\|_2^4\|_{p/2}^{1/2} + pq. \\ 533 \end{aligned} \quad (4)$$

534 Here we used the maximum trick from the next section to bound the max term.

535 A short aside: It would be sweet to split

$$\begin{aligned} 536 \quad \|\sum_k \|U_{(k)}\|_2^4\|_{p/2} &\leq \|\max_k \|U_{(k)}\|_2^2\|_p \sum_k \|U_{(k)}\|_2^2\|_{p/2} \\ 537 \quad &\leq \|\max_k \|U_{(k)}\|_2^2\|_p \|\sum_k \|U_{(k)}\|_2^2\|_p \\ 538 \quad &\leq p \|\|U\|_F^2\|_p, \\ 539 \end{aligned}$$

540 but unfortunately the second factor is $\sqrt{pk} + p$, which means we end up with $p(pk)^{1/4}$
541 term in (4), which is too much. ($p^{3/4}k^{1/4}$ would have been tolerable.) We'll show how
542 to shave this factor p on the \sqrt{pk} term.

¹ $\|\|U\|_F^2\|_p = \|\sum_j M_j X^T X M_j^T\|_p \leq \sqrt{pk}\|X\|_F^2 + p\|X\|^2$.

XX:16 High Probability Tensor Sketch

543 We instead have to work directly on the fourth powers. We use triangle and Bernstein

$$\begin{aligned}
 544 \quad \left\| \sum_k \|U_{(k)}\|_2^4 \right\|_{p/2} &\leq \left\| \sum_k \|U_{(k)}\|_2^4 - E \|U_{(k)}\|_2^4 \right\|_{p/2} + \sum_k E \|U_{(k)}\|_2^4 \\
 545 \quad &\leq \sqrt{p} \left(\sum_k (E \|U_{(k)}\|_2^4)^2 \right)^{1/2} + p \max_k \|U_{(k)}\|_2^4 + \sum_k \| \|U_{(k)}\|_2^2 \|_2^2 \\
 546 \quad &\leq \sqrt{p} \left(\sum_k \| \|U_{(k)}\|_2^2 \|_4^4 \right)^{1/2} + p \max_k \|U_{(k)}\|_2^2 + 4k \\
 547 \quad &\leq \sqrt{p} \left(\sum_k 4^4 \right)^{1/2} + p \max_k \| \|U_{(k)}\|_2^2 \|_q^2 + k \\
 548 \quad &\leq \sqrt{pk} + pq^2 + k. \\
 549
 \end{aligned}$$

550 Now plugging into (2) and (4) we get

$$\begin{aligned}
 551 \quad \|(T \bullet M)x\|_2^2/k - \|x\|_2^2 &\leq (\sqrt{p} \sqrt{\sqrt{pk} + pq^2 + k} + pq)/k + \sqrt{p/k} + p/k \\
 552 \quad &\leq \sqrt{p/k} + pq/k. \\
 553
 \end{aligned}$$

554 Finally we can normalize and plug into the power-Markov inequality:

$$555 \quad \Pr[\|(T \bullet M)x\|_2^2/k - \|x\|_2^2 \geq \epsilon] \leq \max(\sqrt{p/k}, pq/k)^p \epsilon^{-p},$$

557 which gives that we must take

$$558 \quad k = \epsilon^{-2} \log 1/\delta + \epsilon^{-1} (\log 1/\delta)^2.$$

560

561 Proofs of the lemmas:

562 ► **Lemma 28** (Max trick). *Let $q = \max(p, \log k)$, then*

$$563 \quad \|\max X_i\|_p \leq e \max \|EX_i\|_q.$$

Proof.

$$\begin{aligned}
 565 \quad \|\max X_i\|_p &\leq \|\max X_i\|_q \\
 566 \quad &= (E \max X_i^q)^{1/q} \\
 567 \quad &\leq \left(\sum EX_i^q \right)^{1/q} \\
 568 \quad &\leq (k \max EX_i^q)^{1/q} \\
 569 \quad &\leq e (\max EX_i^q)^{1/q} \\
 570 \quad &= e \max \|EX_i\|_q. \\
 571
 \end{aligned}$$

572

573 ► **Lemma 29** (p -norm Bernstein). *For independent variables X_i ,*

$$574 \quad \left\| \sum_i X_i - E \sum_i X_i \right\|_p \leq \sqrt{p} \left(\sum_i EX_i^2 \right)^{1/2} + p \max_i \|X_i\|_p.$$

575

Proof.

$$\begin{aligned}
576 \quad & \left\| \sum_i X_i - E \sum_i X_i \right\|_p \leq \left\| \sum_i g_i X_i \right\|_p && \text{(See notes)} \\
577 \quad & \leq \sqrt{p} \left\| \sum_i X_i^2 \right\|_{p/2}^{1/2} \\
578 \quad & \leq \sqrt{p} \left(E \sum_i X_i^2 \right)^{1/2} + \sqrt{p} \left\| \sum_i X_i^2 - E \sum_i X_i^2 \right\|_{p/2}^{1/2} && \text{(Triangle)} \\
579 \quad & \leq \sqrt{p} \sigma + \sqrt{p} \left\| \sum_i X_i^2 - E \sum_i X_i^2 \right\|_{p/2}^{1/2} \\
580 \quad & \leq \sqrt{p} \sigma + \sqrt{p} \left\| \sum_i g_i X_i^2 \right\|_{p/2}^{1/2} \\
581 \quad & \leq \sqrt{p} \sigma + \sqrt{p} \left\| \max_i X_i \sum_i g_i X_i \right\|_{p/2}^{1/2} \\
582 \quad & \leq \sqrt{p} \sigma + \sqrt{p} \left\| \max_i X_i \right\|_p^{1/2} \left\| \sum_i g_i X_i \right\|_p^{1/2}. && \text{(Cauchy)} \\
583 \quad &
\end{aligned}$$

584 Now let $Q = \left\| \sum_i g_i X_i \right\|_p^{1/2}$ and $K = \left\| \max_i X_i \right\|_p$, and we have $Q^2 \leq \sqrt{p} \sigma + \sqrt{pK} Q$. Because
585 it's a quadratic form, Q is upper bounded by the larger root of $Q^2 - \sqrt{pK} Q - \sqrt{p} \sigma$. By
586 calculation, $Q^2 \leq \sqrt{p} \sigma + pK$, which is the theorem. ◀

587 7.2 Proof of polynomial lemma

588 Proof of Lemma 24

589 **Proof.** For each monomial $\alpha x^S y^T$ of P , where S and T are multisets : $[d] \rightarrow \mathbb{N}$, define
590 $[\chi^S y^T]_x$ and $[\chi^S y^T]_y$ be two vectors in \mathbb{R}^ℓ for some ℓ such that

$$591 \quad \langle [\chi^S y^T]_x, [\chi^S y^T]_y \rangle = \sum_{\pi} x^{\pi S} y^{\pi T} = \sum_{\pi} \prod_{i=1}^d x_i^{S_{\pi i}} y_i^{T_{\pi i}}. \quad (5)$$

593 Then $f(x)$ and $g(y)$ are simple the sketched concatenation of $\alpha [\chi^S y^T]_x$ and $[\chi^S y^T]_y$ vectors.
594 (Note that since we get $\epsilon \|f(x)\|_2 \|g(y)\|_2$ error, it doesn't matter where we put alpha, or if
595 we split it between f and g .)

596 We can let $[\chi^0 y^0]_x = [\chi^0 y^0]_y = [1] \in \mathbb{R}^1$ (the single 1 vector) and then define recursively:

$$\begin{aligned}
597 \quad & [\chi^S y^T]_x = x^{\circ S_i} \otimes [\chi^{S \setminus i} y^{T \setminus i}]_x \oplus \bigoplus_{j \in (S \cup T) \setminus i} [\chi^{S \setminus i + \{j : S_i\}} y^{T \setminus i + \{j : T_i\}}]_x \\
598 \quad & [\chi^S y^T]_y = y^{\circ T_i} \otimes [\chi^{S \setminus i} y^{T \setminus i}]_y \oplus \bigoplus_{j \in (S \cup T) \setminus i} [\chi^{S \setminus i + \{j : S_i\}} y^{T \setminus i + \{j : T_i\}}]_y, \\
599 \quad &
\end{aligned}$$

600 where i is any index in $S \cup T$. Here we let $S \setminus i$ be S with i removed, and $S + \{j : S_i\}$ be S
601 with S_i added to S_j . It is clear from Theorem 3 that this construction gives (5).

602 We note that we can compute the norms by $\|[\chi^0 y^0]\|_2^2 = 1$ and

$$603 \quad \left\| [\chi^S y^T]_x \right\|_2^2 = \left\| x^{\circ S_i} \right\|_2^2 \cdot \left\| [\chi^{S \setminus i} y^{T \setminus i}]_x \right\|_2^2 + \sum_{j \in (S \cup T) \setminus i} \left\| [\chi^{S \setminus i + \{j : S_i\}} y^{T \setminus i + \{j : T_i\}}]_x \right\|_2^2$$

605 and equivalently for y . It does however not seem simple to get a closed form in the general
606 case. In the simple case where x and y are $\{0, 1\}^d$ vectors we can however show the simple

XX:18 High Probability Tensor Sketch

607 formula:

$$608 \quad \|\chi^S y^T\|_x^2 = \frac{(\|x\|_2^2 - 1 + |S \cup T|)!}{(\|x\|_2^2 - 1)!} \leq (\|x\|_2^2 - 1 + |S \cup T|)^{|S \cup T|}$$

609

610 and equivalently for y , (Here S and T are normal sets.)

611 Since there are only $4^{|S \cup T|}$ many states of $[\chi^S y^T]$ the running time is only that many
612 sketching operations.

613 ◀

614 7.3 Better Approximate Matrix Multiplication

615 ▶ **Lemma 30.** For any $x, y \in \mathbb{R}^d$, if S has the (ϵ, δ) -JL-moment-property, $(\|Sx\|_2 - \|x\|_2)_p \leq$
616 $\epsilon \delta^{1/p} \|x\|_2$, then also

$$617 \quad \|(Sx)^T(Sy) - x^T y\|_p \leq \epsilon \delta^{1/p} \|x\|_2 \|y\|_2$$

618

619 **Proof.** We can assume by linearity of the norms that $\|x\|_2 = \|y\|_2 = 1$. We then use that
620 $\|x - y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 - 2x^T y$ and $\|x + y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 + 2x^T y$.

$$621 \quad \begin{aligned} \|(Sx)^T(Sy) - x^T y\|_p &= \| \|Sx + Sy\|_2^2 - \|x + y\|_2^2 - \|Sx - Sy\|_2^2 + \|x - y\|_2^2 \|_p / 4 \\ 622 \quad &\leq (\| \|S(x + y)\|_2^2 - \|x + y\|_2^2 \|_p + \| \|S(x - y)\|_2^2 - \|x - y\|_2^2 \|_p) / 4 \\ 623 \quad &\leq \epsilon \delta^{1/p} (\|x + y\|_2^2 + \|x - y\|_2^2) / 4 \quad (\text{JL property}) \\ 624 \quad &= \epsilon \delta^{1/p} (\|x\|_2^2 + \|y\|_2^2) / 2 \\ 625 \quad &\leq \epsilon \delta^{1/p}. \end{aligned}$$

626

627 Combined with the argument in [28] this gives that the JL-moment-property implies Ap-
628 proximate Matrix Multiplication without a factor 3 on ϵ . ◀