

Oblivious Sketching of High-Degree Polynomial Kernels*

Thomas D. Ahle ITU and BARC thdy@itu.dk	Michael Kapralov EPFL michael.kapralov@epfl.ch	Jakob B. T. Knudsen U. Copenhagen and BARC jagn@di.ku.dk	
Rasmus Pagh ITU and BARC pagh@itu.dk	Ameya Velingker Google Research ameyav@google.com	David P. Woodruff CMU dwoodruf@cs.cmu.edu	Amir Zandieh EPFL amir.zandieh@epfl.ch

April 21, 2020

Abstract

Kernel methods are fundamental tools in machine learning that allow detection of non-linear dependencies between data without explicitly constructing feature vectors in high dimensional spaces. A major disadvantage of kernel methods is their poor scalability: primitives such as kernel PCA or kernel ridge regression generally take prohibitively large quadratic space and (at least) quadratic time, as kernel matrices are usually dense. Some methods for speeding up kernel linear algebra are known, but they all invariably take time exponential in either the dimension of the input point set (e.g., fast multipole methods suffer from the *curse of dimensionality*) or in the degree of the kernel function.

Oblivious sketching has emerged as a powerful approach to speeding up numerical linear algebra over the past decade, but our understanding of oblivious sketching solutions for kernel matrices has remained quite limited, suffering from the aforementioned exponential dependence on input parameters. Our main contribution is a general method for applying sketching solutions developed in numerical linear algebra over the past decade to a tensoring of data points without forming the tensoring explicitly. This leads to the first oblivious sketch for the polynomial kernel with a target dimension that is only polynomially dependent on the degree of the kernel function, as well as the first oblivious sketch for the Gaussian kernel on bounded datasets that does not suffer from an exponential dependence on the dimensionality of input data points.

*This paper is a merged version of the work of Ahle and Knudsen [AK19] and Kapralov, Pagh, Velingker, Woodruff and Zandieh [KPV⁺19].

Contents

1	Introduction	3
1.1	Our Contributions	4
1.2	Technical Overview	8
1.3	Related Work	11
1.4	Organization	13
2	Preliminaries	13
3	Construction of the Sketch	14
4	Linear Dependence on the Tensoring Degree p	18
4.1	Second Moment of Π^q (analysis for T_{base} : CountSketch and S_{base} : TensorSketch) . .	24
4.2	Higher Moments of Π^q (analysis for T_{base} : OSNAP and S_{base} : TensorSRHT)	26
5	Linear Dependence on the Statistical Dimension s_λ	34
5.1	Matrix Concentration Tools	34
5.2	Spectral Property of the sketch Π^q	36
5.3	Spectral Property of Identity \times TensorSRHT	40
5.4	Spectral property of Identity \times OSNAP	43
5.5	High Probability OSE with linear dependence on s_λ	47
6	Oblivious Subspace Embedding for the Gaussian Kernel	48
A	Direct Lower and Upper Bounds	54
A.1	Lower Bound for Sub-Gaussians	55
A.2	Upper bound for Sub-Gaussians	57
A.3	Lower Bound for TensorSketch	58

1 Introduction

Data dimensionality reduction, or *sketching*, is a common technique for quickly reducing the size of a large-scale optimization problem while approximately preserving the solution space, thus allowing one to instead solve a much smaller optimization problem, typically in a smaller amount of time. This technique has led to near-optimal algorithms for a number of fundamental problems in numerical linear algebra and machine learning, such as least squares regression, low rank approximation, canonical correlation analysis, and robust variants of these problems. In a typical instance of such a problem, one is given a large matrix $X \in \mathbb{R}^{d \times n}$ as input, and one wishes to choose a random map Π from a certain family of random maps and replace X with ΠX . As Π typically has many fewer rows than columns, ΠX compresses the original matrix X , which allows one to perform the original optimization problem on the much smaller matrix ΠX . For a survey of such techniques, we refer the reader to the survey by Woodruff [Woo14].

A key challenge in this area is to extend sketching techniques to kernel-variants of the above linear algebra problems. Suppose each column of X corresponds to an example while each of the d rows corresponds to a feature. Then these algorithms require an explicit representation of X to be made available to the algorithm. This is unsatisfactory in many machine learning applications, since typically the actual learning is performed in a much higher (possibly infinite) dimensional feature space, by first mapping each column of X to a much higher dimensional space. Fortunately, due to the kernel trick, one need not ever perform this mapping explicitly; indeed, if the optimization problem at hand only depends on inner product information between the input points, then the kernel trick allows one to quickly compute the inner products of the high dimensional transformations of the input points, without ever explicitly computing the transformation itself. However, evaluating the kernel function easily becomes a bottleneck in algorithms that rely on the kernel trick because it typically takes $O(d)$ time to evaluate the kernel function for d dimensional datasets. There are a number of recent works which try to improve the running times of kernel methods; we refer the reader to the recent work of [MM17] and the references therein. A natural question is whether it is possible to instead apply sketching techniques on the high-dimensional feature space without ever computing the high-dimensional mapping.

For the important case of *polynomial kernel*, such sketching techniques are known to be possible¹. This was originally shown by Pham and Pagh in the context of kernel support vector machines [PP13], using the *TensorSketch* technique for compressed matrix multiplication due to Pagh [Pag13]. This was later extended in [ANW14] to a wide array of kernel problems in linear algebra, including principal component analysis, principal component regression, and canonical correlation analysis.

The running times of the algorithms above, while nearly linear in the number of non-zero entries of the input matrix X , depend *exponentially* on the degree q of the polynomial kernel. For example, suppose one wishes to do low rank approximation on A , the matrix obtained by replacing each column of X with its kernel-transformed version. One would like to express $A \approx UV$, where $U \in \mathbb{R}^{d^p \times k}$ and $V \in \mathbb{R}^{k \times n}$. Writing down U explicitly is problematic, since the columns belong to the much higher d^p -dimensional space. Instead, one can express UV implicitly via column subset selection, by expressing it as a AZZ^T and then outputting Z . Here Z is an $n \times k$ matrix. In [ANW14], an algorithm running in $\text{nnz}(X) + (n+d)\text{poly}(3^p, k, 1/\epsilon)$ time was given for outputting such Z with the guarantee that $\|A - AZZ^T\|_F^2 \leq (1+\epsilon)\|A - A_k\|_F^2$ with constant probability, where A_k is the best rank- k approximation to A . Algorithms with similar running times were proposed for principal component regression and canonical correlation analysis. The main message here is

¹The lifting function corresponding to the polynomial kernel maps $x \in \mathbb{R}^d$ to $\phi(x) \in \mathbb{R}^{d^p}$, where $\phi(x)_{i_1, i_2, \dots, i_p} = x_{i_1} x_{i_2} \cdots x_{i_p}$, for $i_1, i_2, \dots, i_p \in \{1, 2, \dots, d\}$

that all analyses of all existing sketches require the sketch Π to have at least 3^p rows in order to guarantee their correctness. Moreover, the existing sketches work with constant probability only and no high probability result was known for the polynomial kernel.

The main drawback with previous work on applying dimensionality reduction for the polynomial kernel is the exponential dependence on p in the sketching dimension and consequently in the running time. Ideally, one would like a polynomial dependence. This is especially useful for the application of approximating the Gaussian kernel by a sum of polynomial kernels of various degrees, for which large values of p , e.g., $p = \text{poly}(\log n)$ are used [CKS11]. This raises the main question of our work:

Is it possible to design a data oblivious sketch with a sketching dimension (and, hence, running time) that is not exponential in p for the above applications in the context of the polynomial kernel?

While we answer the above question, we also study it in a more general context, namely, that of regularization. In many machine learning problems, it is crucial to regularize so as to prevent overfitting or ill-posed problems. Sketching and related sampling-based techniques have also been extensively applied in this setting. For a small sample of such work see [RR07, AM15, PW15, MM17, ACW17b, ACW17a, AKM⁺17, AKM⁺18a]. As an example application, in ordinary least squares regression one is given a $d \times n$ matrix A , and a $d \times 1$ vector b , and one seeks to find a $y \in \mathbb{R}^n$ so as to minimize $\|Ay - b\|_2^2$. In ridge regression, we instead seek a y so as to minimize $\|Ay - b\|_2^2 + \lambda\|y\|_2^2$, for a parameter $\lambda > 0$. Intuitively, if λ is much larger than the operator norm $\|A\|_2$ of A , then a good solution is obtained simply by setting $y = 0^d$. On the other hand, if $\lambda = 0$, the problem just becomes an ordinary least squares regression. In general, the *statistical dimension* (or *effective degrees of freedom*), s_λ , captures this tradeoff, and is defined as $\sum_{i=1}^d \frac{\lambda_i(A^\top A)}{\lambda_i(A^\top A) + \lambda}$, where $\lambda_i(A^\top A)$ is the i -th eigenvalue of $A^\top A$. Note that the statistical dimension is always at most $\min(n, d)$, but in fact can be much smaller. A key example of its power is that for ridge regression, it is known [ACW17b] that if one chooses a random Gaussian matrix Π with $O(s_\lambda/\epsilon)$ rows, and if y is the minimizer to $\|\Pi Ay - \Pi b\|_2^2 + \lambda\|y\|_2^2$, then $\|Ay - b\|_2^2 + \lambda\|y\|_2^2 \leq (1 + \epsilon) \min_{y'} (\|Ay' - b\|_2^2 + \lambda\|y'\|_2^2)$. Note that for ordinary regression ($\lambda = 0$) one would need that Π has $\Omega(\text{rank}(A)/\epsilon)$ rows [CW09]. Another drawback of existing sketches for the polynomial kernel is that their running time and target dimension depend at least quadratically on s_λ and no result is known with linear dependence on s_λ , which would be optimal. We also ask if the exponential dependence on p is avoidable in the *regularized* setting:

Is it possible to obtain sketching dimension bounds and running times that are not exponential in p in the context of regularization? Moreover, is it possible to obtain a running time that depends only linearly on s_λ ?

1.1 Our Contributions

In this paper, we answer the above questions in the affirmative. In other words, for each of the aforementioned applications, our algorithm depends only *polynomially* on p . We state these applications as corollaries of our main results, which concern approximate matrix product and subspace embeddings. In particular, we devise a new distribution on oblivious linear maps $\Pi \in \mathbb{R}^{m \times d^p}$ (i.e., a randomized family of maps that does not depend on the dataset X), so that for any fixed $X \in \mathbb{R}^{d \times n}$, it satisfies the approximate matrix product and subspace embedding properties. These are the key properties needed for kernel low rank approximation. We remark that our

data oblivious sketching is greatly advantageous to data dependent methods because it results in a one-round distributed protocol for kernel low rank approximation [KVV14].

We show that our oblivious linear map $\Pi \in \mathbb{R}^{m \times d^p}$ has the following key properties:

Oblivious Subspace Embeddings (OSEs). Given $\varepsilon > 0$ and an n -dimensional subspace $E \subseteq \mathbb{R}^d$, we say that $\Pi \in \mathbb{R}^{m \times d}$ is an ε -subspace embedding for E if $(1 - \varepsilon)\|x\|_2 \leq \|\Pi x\|_2 \leq (1 + \varepsilon)\|x\|_2$ for all $x \in E$. In this paper we focus on Oblivious Subspace Embeddings in the regularized setting. In order to define a (regularized) Oblivious Subspace Embedding, we need to introduce the notion of *statistical dimension*, which is defined as follows:

Definition 1 (Statistical Dimension). Given $\lambda \geq 0$, for every positive semidefinite matrix $K \in \mathbb{R}^{n \times n}$, we define the λ -statistical dimension of K to be

$$s_\lambda(K) := \text{tr}(K(K + \lambda I_n)^{-1}).$$

Now, we can define the notion of an oblivious subspace embedding (OSE):

Definition 2 (Oblivious Subspace Embedding (OSE)). Given $\varepsilon, \delta, \mu > 0$ and integers $d, n \geq 1$, an $(\varepsilon, \delta, \mu, d, n)$ -Oblivious Subspace Embedding (OSE) is a distribution \mathcal{D} over $m \times d$ matrices (for arbitrary m) such that for every $\lambda \geq 0$, every $A \in \mathbb{R}^{d \times n}$ with λ -statistical dimension $s_\lambda(A^\top A) \leq \mu$, the following holds,²

$$\Pr_{\Pi \sim \mathcal{D}} \left[(1 - \varepsilon)(A^\top A + \lambda I_n) \preceq (\Pi A)^\top \Pi A + \lambda I_n \preceq (1 + \varepsilon)(A^\top A + \lambda I_n) \right] \geq 1 - \delta. \quad (1)$$

The goal is to have the target dimension m small so that Π provides dimensionality reduction. If we consider the non-oblivious setting where we allow the sketch matrix Π to depend on A , then by leverage score sampling we can achieve a target dimension of $m \approx s_\lambda(A^\top A)$, which is essentially optimal [AKM⁺18b]. But as we discussed the importance of oblivious embeddings, the ultimate goal is to get an oblivious subspace embedding with target dimension of $m \approx s_\lambda(A^\top A)$.

Approximate Matrix Product. We formally define this property in the following definition.

Definition 3 (Approximate Matrix Product). Given $\varepsilon, \delta > 0$, we say that a distribution \mathcal{D} over $m \times d$ matrices has the (ε, δ) -approximate matrix product property if for every $C, D \in \mathbb{R}^{d \times n}$,

$$\Pr_{\Pi \sim \mathcal{D}} \left[\|C^\top \Pi^\top \Pi D - C^\top D\|_F \leq \varepsilon \|C\|_F \|D\|_F \right] \geq 1 - \delta.$$

Our main theorems, which provide the aforementioned guarantees, are as follows,³

Theorem 1. For every positive integers n, p, d , every $\varepsilon, s_\lambda > 0$, there exists a distribution on linear sketches $\Pi^p \in \mathbb{R}^{m \times d^p}$ such that: **(1)** If $m = \Omega(p s_\lambda^2 \varepsilon^{-2})$, then Π^p is an $(\varepsilon, 1/10, s_\lambda, d^p, n)$ -oblivious subspace embedding as in Definition 2. **(2)** If $m = \Omega(p \varepsilon^{-2})$, then Π^p has the $(\varepsilon, 1/10)$ -approximate matrix product property as in Definition 3.

Moreover, for any $X \in \mathbb{R}^{d \times n}$, if $A \in \mathbb{R}^{d^p \times n}$ is the matrix whose columns are obtained by the p -fold self-tensoring of each column of X then the matrix $\Pi^p A$ can be computed using Algorithm 1 in time $\tilde{O}(pnm + p \text{nnz}(X))$.

²For symmetric matrices K and K' , the spectral inequality relation $K \preceq K'$ holds if and only if $x^\top K x \leq x^\top K' x$ for all vectors x

³Throughout this paper, the notations $\tilde{O}, \tilde{\Omega}, \tilde{\Theta}$ suppress poly($\log(nd/\varepsilon)$) factors.

Theorem 2. For every positive integers n, p, d , every $\varepsilon, s_\lambda > 0$, there exists a distribution on linear sketches $\Pi^p \in \mathbb{R}^{m \times d^p}$ such that: **(1)** If $m = \tilde{\Omega}(ps_\lambda^2 \varepsilon^{-2})$, then Π^p is an $(\varepsilon, 1/\text{poly}(n), s_\lambda, d^p, n)$ -oblivious subspace embedding (Definition 2). **(2)** If $m = \tilde{\Omega}(p\varepsilon^{-2})$, then Π^p has the $(\varepsilon, 1/\text{poly}(n))$ -approximate matrix product property (Definition 3).

Moreover, in the setting of **(1)**, for any $X \in \mathbb{R}^{d \times n}$, if $A \in \mathbb{R}^{d^p \times n}$ is the matrix whose columns are obtained by a p -fold self-tensoring of each column of X , then the matrix $\Pi^p A$ can be computed using Algorithm 1 in time $\tilde{O}(pnm + p^{3/2}s_\lambda \varepsilon^{-1} \text{nnz}(X))$.

Theorem 3. For every positive integers p, d, n , every $\varepsilon, s_\lambda > 0$, there exists a distribution on linear sketches $\Pi^p \in \mathbb{R}^{m \times d^p}$ which is an $(\varepsilon, 1/\text{poly}(n), s_\lambda, d^p, n)$ -oblivious subspace embedding as in Definition 2, provided that the integer m satisfies $m = \tilde{\Omega}(p^4 s_\lambda / \varepsilon^2)$.

Moreover, for any $X \in \mathbb{R}^{d \times n}$, if $A \in \mathbb{R}^{d^p \times n}$ is the matrix whose columns are obtained by a p -fold self-tensoring of each column of X then the matrix $\Pi^p A$ can be computed using Algorithm 1 in time $\tilde{O}(pnm + p^5 \varepsilon^{-2} \text{nnz}(X))$.

We can immediately apply these theorems to *kernel ridge regression* with respect to the polynomial kernel of degree p . In this problem, we are given a regularization parameter $\lambda > 0$, a $d \times n$ matrix X , and vector $b \in \mathbb{R}^n$ and would like to find a $y \in \mathbb{R}^n$ so as to minimize $\|A^\top Ay - b\|_2^2 + \lambda \|Ay\|_2^2$, where $A \in \mathbb{R}^{d^p \times n}$ is the matrix obtained from X by applying the self tensoring of degree p to each column. To solve this problem via sketching, we choose a random matrix Π^p according to the theorems above and compute $\Pi^p A$. We then solve the sketched ridge regression problem which seeks to minimize $\|(\Pi^p A)^\top \Pi^p Ax - b\|_2^2 + \lambda \|\Pi^p Ax\|_2^2$ over x . By the above theorems, we have $\|(\Pi^p A)^\top \Pi^p Ax - b\|_2^2 + \lambda \|\Pi^p Ax\|_2^2 = (1 \pm \varepsilon) \left(\|A^\top Ax - b\|_2^2 + \lambda \|Ax\|_2^2 \right)$ simultaneously for all $x \in \mathbb{R}^n$; thus, solving the sketched ridge regression problem gives a $(1 \pm \varepsilon)$ -approximation to the original problem. If we apply Theorem 1, then the number of rows of Π^p needed to ensure success with probability $9/10$ is $\Theta(ps_\lambda^2 \varepsilon^{-2})$. The running time to compute $\Pi^p A$ is $O(p^2 s_\lambda^2 \varepsilon^{-2} n + p \text{nnz}(X))$, after which a ridge regression problem can be solved in $O(ns_\lambda^4 / \varepsilon^4)$ time via an exact closed-form solution for ridge regression. An alternative approach to obtaining a very high-accuracy approximation is to use the sketched kernel as a preconditioner to solve the original ridge regression problem, which improves the dependence on ε to $\log(1/\varepsilon)$ [ACW17a]. To obtain a higher probability of success, we can instead apply Theorem 3, which would allow us to compute the sketched matrix $\Pi^p A$ in $\tilde{O}(p^5 s_\lambda \varepsilon^{-2} n + p^5 \varepsilon^{-2} \text{nnz}(X))$ time. This is the first sketch to achieve the optimal dependence on s_λ for the polynomial kernel, after which we can now solve the ridge regression problem in $\tilde{O}(ns_\lambda^2 \text{poly}(p, \varepsilon^{-1}))$ time. Importantly, both running times are polynomial in p , whereas all previously known methods incurred running times that were exponential in p .

Although there has been much work on sketching methods for kernel approximation which nearly achieve the optimal target dimension $m \approx s_\lambda$, such as Nystrom sampling [MM17], all known methods are data-dependent unless strong conditions are assumed about the kernel matrix (small condition number or incoherence). Data oblivious methods provide nice advantages, such as one-round distributed protocols and single-pass streaming algorithms. However, for kernel methods they are poorly understood and previously had worse theoretical guarantees than data-dependent methods. Furthermore, note that the Nystrom method requires to sample at least $m = \Omega(s_\lambda)$ landmarks to satisfy the subspace embedding property even given an oracle access to the exact leverage scores distribution. This results in a runtime of $\Omega(s_\lambda^2 d + s_\lambda \text{nnz}(X))$. Whereas our method achieves a target dimension that nearly matches the best dimension possible with data-dependent Nystrom method and with strictly better running time of $\tilde{O}(ns_\lambda + \text{nnz}(X))$ (assuming $p = \text{poly}(\log n)$).

Therefore, for a large range of parameter our sketch runs in input sparsity time whereas the Nystrom methods are slower by an s_λ factor in the best case.

Application: Polynomial Kernel Rank- k Approximation. Approximate matrix product and subspace embedding are key properties for sketch matrices which imply efficient algorithms for rank- k kernel approximation [ANW14]. The following corollary of Theorem 1 immediately follows from Theorem 6 of [ANW14].

Corollary 4 (Rank- k Approximation). *For every positive integers k, n, p, d , every $\epsilon > 0$, any $X \in \mathbb{R}^{d \times n}$, if $A \in \mathbb{R}^{d^p \times n}$ is the matrix whose columns are obtained by the p -fold self-tensoring of each column of X then there exists an algorithm which finds an $n \times k$ matrix V in time $O(p \text{nnz}(X) + \text{poly}(k, p, \epsilon^{-1}))$ such that with probability $9/10$,*

$$\|A - AVV^\top\|_F^2 \leq (1 + \epsilon) \min_{\substack{U \in \mathbb{R}^{d^p \times n} \\ \text{rank}(U)=k}} \|A - U\|_F^2.$$

Note that this runtime improves the runtime of [ANW14] by exponential factors in the polynomial kernel's degree p .

Additional Applications. Our results also imply improved bounds for each of the applications in [ANW14], including canonical correlation analysis (CCA), and principal component regression (PCR). Importantly, we obtain the first sketching-based solutions for these problems with running time polynomial rather than exponential in p .

Oblivious Subspace Embedding for the Gaussian Kernel. One very important implication of our result is Oblivious Subspace Embedding of the Gaussian kernel. Most work in this area is related to the Random Fourier Features method [RR07]. It was shown in [AKM⁺17] that one requires $\Omega(n)$ samples of the standard Random Fourier Features to obtain a subspace embedding for the Gaussian kernel, while a modified distribution for sampling frequencies yields provably better performance. The target dimension of our proposed sketch for the Gaussian kernel strictly improves upon the result of [AKM⁺17], which has an exponential dependence on the dimension d . We for the first time, embed the Gaussian kernel with a target dimension which has a linear dependence on the statistical dimension of the kernel and is not exponential in the dimensionality of the data-point.

Theorem 5. *For every $r > 0$, every positive integers n, d , and every $X \in \mathbb{R}^{d \times n}$ such that $\|x_i\|_2 \leq r$ for all $i \in [n]$, where x_i is the i^{th} column of X , suppose $G \in \mathbb{R}^{n \times n}$ is the Gaussian kernel matrix – i.e., $G_{j,k} = e^{-\|x_j - x_k\|_2^2/2}$ for all $j, k \in [n]$. There exists an algorithm which computes $S_g(X) \in \mathbb{R}^{m \times n}$ in time $\tilde{O}(q^6 \epsilon^{-2} n s_\lambda + q^6 \epsilon^{-2} \text{nnz}(X))$ such that for every $\epsilon, \lambda > 0$,*

$$\Pr_{S_g} \left[(1 - \epsilon)(G + \lambda I_n) \preceq (S_g(X))^\top S_g(X) + \lambda I_n \preceq (1 + \epsilon)(G + \lambda I_n) \right] \geq 1 - 1/\text{poly}(n),$$

where $m = \tilde{\Theta}(q^5 s_\lambda / \epsilon^2)$ and $q = \Theta(r^2 + \log(n/\epsilon\lambda))$ and s_λ is λ -statistical dimension of G as in Definition 1.

We remark that for datasets with radius $r = \text{poly}(\log n)$ even if one has oracle access to the exact leverage scores for Fourier features of Gaussian kernel, in order to get subspace embedding guarantee one needs to use $m = \Omega(s_\lambda)$ features which requires $\Omega(s_\lambda \text{nnz}(X))$ operations to compute. Whereas our result of Theorem 5 runs in time $\tilde{O}(n s_\lambda + \text{nnz}(X))$. Therefore, for a large range of parameters our Gaussian sketch runs in input sparsity time whereas the Fourier features method is at best slower by an s_λ factor.

1.2 Technical Overview

Our goal is to design a sketching matrix Π^p that satisfies the oblivious subspace embedding property with an optimal embedding dimension and which can be efficiently applied to vectors of the form $x^{\otimes p} \in \mathbb{R}^{d^p}$. We start by describing some natural approaches to this problem (some of which have been used before), and show why they incur an exponential loss in the degree of the polynomial kernel. We then present our sketch and outline our proof of its correctness.

We first discuss two natural approaches to tensoring classical sketches, namely the Johnson-Lindenstrauss transform and the CountSketch. We show that both lead to an exponential dependence of the target dimension on p and then present our new approach.

Tensoring the Johnson-Lindenstrauss Transform. Perhaps the most natural approach to designing a sketch Π^p is the idea of tensoring p independent Johnson-Lindenstrauss matrices. Specifically, let m be the target dimension. For every $r = 1, \dots, p$ let $M^{(r)}$ denote an $m \times d$ matrix with iid uniformly random ± 1 entries, and let the sketching matrix $M \in \mathbb{R}^{m \times d^p}$ be

$$M = \frac{1}{\sqrt{m}} M^{(1)} \bullet \dots \bullet M^{(p)},$$

where \bullet stands for the operation of tensoring the rows of matrices $M^{(r)}$ (see Definition 7). This would be a very efficient matrix to apply, since for every $j = 1, \dots, m$ the j -th entry of $Mx^{\otimes p}$ is exactly $\prod_{r=1}^p [M^{(r)}x]_j$, which can be computed in time $O(p \text{nnz}(x))$, giving overall evaluation time $O(pm \text{nnz}(x))$. One would hope that $m = O(\varepsilon^{-2} \log n)$ would suffice to ensure that $\|Mx^{\otimes p}\|_2^2 = (1 \pm \varepsilon) \|x^{\otimes p}\|_2^2$. However, this is not true: we show in Appendix A that one must have $m = \Omega(\varepsilon^{-2} 3^p \log(n)/p + \varepsilon^{-1} (\log(n)/p)^p)$ in order to preserve the norm with high probability. Thus, the dependence on degree p of the polynomial kernel must be exponential. The lower bound is provided by controlling the moments of the sketch M and using Paley-Zygmund inequality. For completeness, we show that the aforementioned bound on the target dimension m is sharp, i.e., necessary and sufficient for obtaining the Johnson-Lindenstrauss property.

Tensoring of CountSketch (TensorSketch). Pagh and Pham [PP13] introduced the following tensorized version of CountSketch. For every $i = 1, \dots, p$ let $h_i : [d] \rightarrow [m]$ denote a random hash function, and $\sigma_i : [d] \rightarrow [m]$ a random sign function. Then let $S : \mathbb{R}^{d^{\otimes p}} \rightarrow \mathbb{R}^m$ be defined by

$$S_{r,(j_1, \dots, j_p)} := \sigma(i_1) \cdots \sigma(i_p) \mathbf{1}[h_1(i_1) + \dots + h_p(i_p) = r]$$

for $r = 1, \dots, m$. For every $x \in \mathbb{R}^d$ one can compute $Sx^{\otimes p}$ in time $O(pm \log m + p \text{nnz}(x))$. Since the time to apply the sketch only depends linearly on the dimension p (due to the Fast Fourier Transform) one might hope that the dependence of the sketching dimension on p is polynomial. However, this turns out to not be the case: the argument in [ANW14] implies that $m = \tilde{O}(3^p s_\lambda^2)$ suffices to construct a subspace embedding for a matrix with regularization λ and statistical dimension s_λ , and we show in Appendix A.3 that exponential dependence on p is necessary.

Our Approach: Recursive Tensoring. The initial idea behind our sketch is as follows. To apply our sketch Π^p to $x^{\otimes p}$, for $x \in \mathbb{R}^d$, we first compute the sketches T_1x, T_2x, \dots, T_px for independent sketching matrices $T_1, \dots, T_p \sim T_{\text{base}}$ – see the leaves of the sketching tree in Fig. 1. Note that we choose these sketches as CountSketch [CCFC02] or OSNAP [NN13] to ensure that the leaf

⁴Tensor product of x with itself p times.

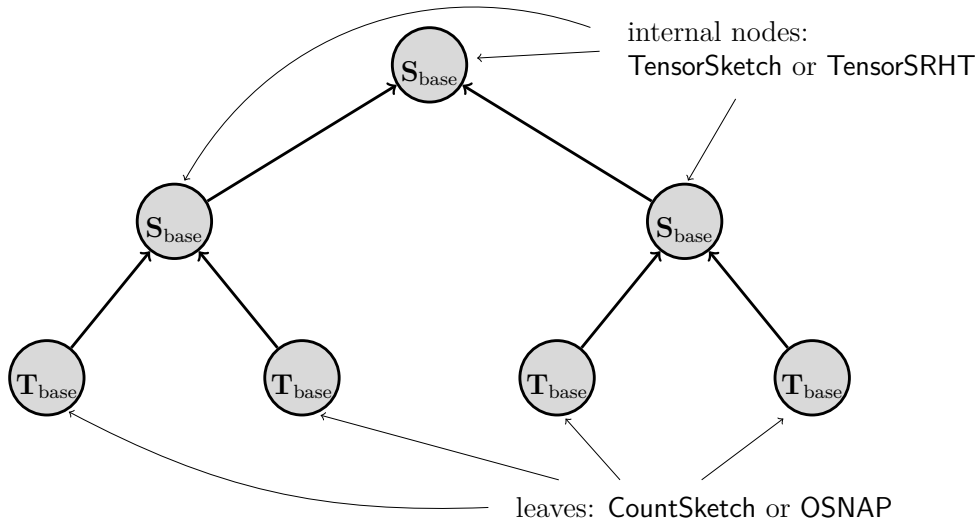


Figure 1: S_{base} is chosen from the family of sketches which support fast matrix-vector product for tensor inputs such as TensorSketch and TensorSRHT. The T_{base} is chosen from the family of sketches which operate in input sparsity time such as CountSketch and OSNAP.

sketches can be applied in time proportional to the number of nonzeros in the input data (in the case of OSNAP this is true up to polylogarithmic factors).

Each of these is a standard sketching matrix mapping d -dimensional vectors to m -dimensional vectors for some common value of m . We refer the reader to the survey [Woo14]. The next idea is to choose new sketching matrices $S_1, S_2, \dots, S_{p/2} \sim S_{\text{base}}$, mapping m^2 -dimensional vectors to m -dimensional vectors and apply S_1 to $(T_1x) \otimes (T_2x)$, as well as apply S_2 to $(T_3x) \otimes (T_4x)$, and so on, applying $S_{p/2}$ to $(T_{p-1}x) \otimes (T_px)$. These sketches are denoted by S_{base} – see internal nodes of the sketching tree in Fig. 1. We note that in order to ensure efficiency of our construction (in particular, running time that depends only linearly on the statistical dimension s_λ) we must choose S_{base} as a sketch that can be computed on tensored data without explicitly constructing the actual tensored input, i.e., S_{base} supports fast matrix vector product on tensor product of vectors. We use either TensorSketch (for results that work with constant probability) and a new variant of the Subsampled Randomized Hadamard Transform SRHT which supports fast multiplication for the tensoring of two vectors (for high probability bounds) – we call the last sketch TensorSRHT.

At this point we have reduced our number of input vectors from p to $p/2$, and the dimension is m , which will turn out to be roughly s_λ . We have made progress, as we now have fewer vectors each in roughly the same dimension we started with. After $\log_2 p$ levels in the tree we are left with a single output vector.

Intuitively, the reason that this construction avoids an exponential dependence on p is that at every level in the tree we use target dimension m larger than the statistical dimension of our matrix by a factor polynomial in p . This ensures that the accumulation of error is limited, as the total number of nodes in the tree is $O(p)$. This is in contrast to the direct approaches discussed above, which use a rather direct tensoring of classical sketches, thereby incurring an exponential dependence on p due to dependencies that arise.

Showing Our Sketch is a Subspace Embedding. In order to show that our recursive sketch is a subspace embedding, we need to argue it preserves norms of arbitrary vectors in \mathbb{R}^{d^p} , not only

vectors of the form $x^{\otimes p}$, i.e., p -fold self-tensoring of d -dimensional vectors⁵. Indeed, all known methods for showing the subspace embedding property (see [Woo14] for a survey) at the very least argue that the norms of each of the columns of an orthonormal basis for the subspace in question are preserved. While our subspace may be formed by the span of vectors which are tensor products of p d -dimensional vectors, we are not guaranteed that there is an orthonormal basis of this form. Thus, we first observe that our mapping is indeed linear over \mathbb{R}^{d^p} , making it well-defined on the elements of any basis for our subspace, and hence our task essentially reduces to proving that our mapping preserves norms of arbitrary vectors in \mathbb{R}^{d^p} .

We present two approaches to analyzing our construction. One is based on the idea of propagating moment bounds through the sketching tree, and results in a nearly linear dependence of the sketching dimension m on the degree p of the polynomial kernel, at the expense of a quadratic dependence on the statistical dimension s_λ . This approach is presented in Section 4. The other approach achieves the (optimal) linear dependence on s_λ , albeit at the expense of a worse polynomial dependence on p . This approach uses sketches that succeed with high probability, and uses matrix concentration bounds.

Propagating moment bounds through the tree – optimizing the dependence on the degree p . We analyze our recursively tensored version of the OSNAP and CountSketch by showing how moment bounds can be propagated through the tree structure of the sketch. This analysis is presented in Section 4, and results in the proof of Theorem 1 as well as the first part of Theorem 3. The analysis obtained this way give particularly sharp dependencies on p and $\log 1/\delta$.

The idea is to consider the unique matrix $M \in \mathbb{R}^{m \times d^p}$ that acts on simple tensors in the way we have described it recursively above. This matrix could in principle be applied to any vector $x \in \mathbb{R}^{d^p}$ (though it would be slow to realise). We can nevertheless show that this matrix has the (ε, δ, t) -JL Moment Property, which is for parameters $\varepsilon, \delta \in [0, 1], t \geq 2$, and every $x \in \mathbb{R}^d$ with $\|x\|_2 = 1$ the statement $\mathbb{E} \left[\left| \|Mx\|_2^2 - 1 \right|^t \right] \leq \varepsilon^t \delta$.

It can be shown that M is built from our various S_{base} and T_{base} matrices using three different operations: multiplication, direct sum, and row-wise tensoring. In other words, it is sufficient to show that if Q and Q' both have the (ε, δ, t) -JL Moment Property, then so does QQ' , $Q \oplus Q'$ and $Q \bullet Q'$. This turns out to hold for $Q \oplus Q'$, but QQ' and $Q \bullet Q'$ are more tricky. (Here \oplus is the direct sum and \bullet is the composition of tensoring the rows. See section 2 on notation.)

For multiplication, a simple union bound allows us to show that $Q^{(1)}Q^{(2)} \dots Q^{(p)}$ has the $(p\varepsilon, p\delta, t)$ -JL Moment Property. This would unfortunately mean a factor of p^2 in the final dimension. The union bound is clearly suboptimal, since implicitly it assumes that all the matrices conspire to either shrink or increase the norm of a vector, while in reality with independent matrices, we should get a random walk on the real line. Using an intricate decoupling argument, we show that this is indeed the case, and that $Q^{(1)}Q^{(2)} \dots Q^{(p)}$ has the $(\sqrt{p}\varepsilon, \delta, t)$ -JL Moment Property, saving a factor of p in the output dimension.

Finally we need to analyze $Q \bullet Q'$. Here it is easy to show that the JL Moment Property doesn't in general propagate to $Q \bullet Q'$ (consider e.g. Q being constant 0 on its first $m/2$ rows and Q' having 0 on its $m/2$ last rows.) For most known constructions of JL matrices it does however turn out that $Q \bullet Q'$ behaves well. In particular we show this for matrices with independent sub-Gaussian entries (appendix A.2), and for the so-called Fast Johnson Lindenstrauss construction [AC06] (Lemma 21). The main tool here is a higher order version of the classical Khintchine

⁵ $x^{\otimes p}$ denotes $\underbrace{x \otimes x \dots \otimes x}_p$, the p -fold self-tensoring of x .

inequality [HM07] which bounds the moments $\mathbb{E}\left[\langle \sigma^{(1)} \otimes \sigma^{(2)} \otimes \dots \otimes \sigma^{(p)}, x \rangle^t\right]$ when $\sigma^{(1)}, \dots, \sigma^{(p)}$ are independent sub-Gaussian vectors (Lemma 19).

Optimizing the dependence on s_λ . Our proof of Theorem 3 relies on instantiating our framework with OSNAP at the leaves of the tree (T_{base}) and a novel version of the SRHT that we refer to as TensorSRHT at the internal nodes of the tree. We outline the analysis here. In order to show that our sketch preserves norms, let y be an arbitrary vector in \mathbb{R}^{d^p} . Then in the bottom level of the tree, we can view our sketch as $T_1 \times T_2 \times \dots \times T_p$, where \times for denotes the tensor product of matrices (see Definition 5). Then, we can reshape y to be a $d^{q-1} \times d$ matrix Y , and the entries of $T_1 \times T_2 \times \dots \times T_p y$ are in bijective correspondence with those of $T_1 \times T_2 \times \dots \times T_{p-1} Y T_p^\top$. By definition of T_p , it preserves the Frobenius norm of Y , and consequently, we can replace Y with $Y T_p^\top$. We next look at $(T_1 \times T_2 \times \dots \times T_{p-2}) Z (I_d \times T_{p-1}^\top)$, where Z is the $d^{p-2} \times d^2$ matrix with entries in bijective correspondence with those of $Y T_p^\top$. Then we know that T_{p-1} preserves the Frobenius norm of Z . Iterating in this fashion, this means the first layer of our tree preserves the norm of y , provided we union bound over $O(p)$ events that a sketch preserves a norm of an intermediate matrix. The core of the analysis consists of applying spectral concentration bounds based analysis to sketches that act on blocks of the input vector in a correlated fashion. We give the details in Section 5.

Sketching the Gaussian kernel. Our techniques yield the first oblivious sketching method for the Gaussian kernel with target dimension that does not depend exponentially on the dimensionality of the input data points. The main idea is to Taylor expand the Gaussian function and apply our sketch for the polynomial kernel to the elements of the expansion. It is crucial here that the target dimension of our sketch for the polynomial kernel depends only polynomially on the degree, as otherwise we would not be able to truncate the Taylor expansion sufficiently far in the tail (the number of terms in the Taylor expansion depends on the radius of the dataset and depends logarithmically on the regularization parameter). Overall, our Gaussian kernel sketch has optimal target dimension up to polynomial factors in the radius dataset and logarithmic factors in the dataset size. Moreover, it is the first subspace embedding of Gaussian kernel which runs in input sparsity time $\tilde{O}(\text{nnz}(X))$ for datasets with polylogarithmic radius. The result is summarized in Theorem 5, and the analysis is presented in Section 6.

1.3 Related Work

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

Johnson-Lindenstrauss Transform A cornerstone result in the field of subspace embeddings is the Johnson-Lindenstrauss lemma [JLS86]: “For all $\varepsilon \in [0, 1]$, integers $n, d \geq 1$, and $X \subseteq \mathbb{R}^d$ with $|X| = n$ there exists $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$ with $m = O(\varepsilon^{-2} \log(n))$, such that $(1 - \varepsilon)\|x - y\|_2 \leq \|f(x) - f(y)\|_2 \leq (1 + \varepsilon)\|x - y\|_2$ for every $x, y \in X$.”

It has been shown in [CW13, CNW16a] there exists a constant C , so that, for any r -dimensional subspace $U \subseteq \mathbb{R}^d$, there exists a subset $X \subseteq U$ with $|X| \leq C^r$, such that $\max_{x \in U} \left| \|f(x)\|_2^2 - \|x\|_2^2 \right| \leq O(\max_{x \in X} \left| \|f(x)\|_2^2 - \|x\|_2^2 \right|)$. So the Johnson-Lindenstrauss Lemma implies that there exists a subspace embedding with $m = O(\varepsilon^{-2} r)$.

It is not enough to know that the subspace embedding exists, we also need the to find the dimension-reducing map f , and we want the map f to be applied to the data quickly. Achlioptas showed that if $\Pi \in \mathbb{R}^{m \times d}$ is random matrix with i.i.d. entries where $\Pi_{i,j} = 0$ with probability $2/3$, and otherwise $\Pi_{i,j}$ is uniform in $\{-1, 1\}$, and $m = O(\varepsilon^{-2} \log(1/\delta))$, then $\|\Pi x\|_2 = (1 \pm \varepsilon)\|x\|_2$ with probability $1 - \delta$ for any $x \in \mathbb{R}^d$ [Ach03]. This gives a running time of $O(m \text{nnz}(x))$ to sketch a vector $x \in \mathbb{R}^d$. Later, the Fast Johnson Lindenstrauss Transform [AC06], which exploits the Fast Fourier Transform, improved the running time for dense vectors to $O(d \log d + m^3)$. The related Subsampled Randomized Hadamard Transform has been extensively studied [Sar06, DMM06, DMMS11, Tro11, DMMW12, LDFU13], which uses $O(d \log d)$ time but obtains suboptimal dimension $O(\varepsilon^{-2} \log(1/\delta)^2)$, hence it can not use the above argument to get subspace embedding, but it has been proven in [Tro11] that if $m = O(\varepsilon^{-2}(r + \log(1/\delta)^2))$, then one get a subspace embedding.

The above improvements has a running time of $O(d \log d)$, which can be worse than $O(m \text{nnz}(x))$ if $x \in \mathbb{R}^d$ is very sparse. This inspired a line of work trying to obtain sparse Johnson Lindenstrauss transforms [DKS10, KN14, NN13, Coh16]. They obtain a running time of $O(\varepsilon^{-1} \log(1/\delta) \text{nnz}(x))$. In [NN13] they define the ONSAP transform and investigate the trade-off between sparsity and subspace embedding dimension. This was further improved in [Coh16].

In the context of this paper all the above mentioned methods have the same shortcoming, they do not exploit the extra structure of the tensors. The Subsampled Randomized Hadamard Transform have a running time of $\Omega(pd^p \log(p))$ in the model considered in this paper, and the sparse embeddings have a running time of $\Omega(\text{nnz}(x)^p)$. This is clearly unsatisfactory and inspired the TensorSketch [PP13, ANW14], which has a running time of $\Omega(p \text{nnz}(x))$. Unfortunately, they need $m = \Omega(3^p \varepsilon^{-2} \delta^{-1})$ and one of the main contributions of this paper is get rid of the exponential dependence on p .

Approximate Kernel Expansions A classic result by Rahimi and Recht [RR08] shows how to compute an embedding for any shift-invariant kernel function $k(\|x-y\|_2)$ in time $O(dm)$. In [LSS14] this is improved to any kernel on the form $k(\langle x, y \rangle)$ and time $O((m+d) \log d)$, however the method does not handle kernel functions that can't be specified as a function of the inner product, and it doesn't provide subspace embeddings. See also [MM17] for more approaches along the same line. Unfortunately, these methods are unable to operate in input sparsity time and their runtime at best is off by an s_λ factor.

Tensor Sparsification There is also a literature of tensor sparsification based on sampling [NDT15], however unless the vectors tensored are already very smooth (such as ± 1 vectors), the sampling has to be weighted by the data. This means that these methods in aren't applicable in general to the types of problems we consider, where the tensor usually isn't known when the sketching function is sampled.

Hyper-plane rounding An alternative approach is to use hyper-plane rounding to get vectors on the form ± 1 . Let $\rho = \frac{\langle x, y \rangle}{\|x\| \|y\|}$, 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 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 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 $m = O(\rho^{-2} \epsilon^{-2} \log 1/\delta)$ then suffices with high probability. After this we can simply sample from the tensor product using simple sample bounds.

The sign-sketch was first brought into the field of data-analysis by [Cha02] and [Val15] was the first, in our knowledge, to use it with tensoring. The main issue with this approach is that it isn't a linear sketch, which hinders the applications we consider in this paper, such as kernel low rank

approximation, CCA, PCR, and ridge regression. It also takes dm time to calculate Mx and My which is unsatisfactory.

1.4 Organization

In section 2 we introduce basic definitions and notations that will be used throughout the paper. Section 3 introduces our recursive construction of the sketch which is our main technical tool for sketching high degree tensor products. Section 4 analyzes how the moment bounds propagate through our recursive construction thereby proving Theorems 1 and 2 which have linear dependence on the degree q . Section 5 introduces a high probability Oblivious Subspace Embedding with linear dependence on the statistical dimension thereby proving Theorem 3. Finally, section 6 uses the tools that we build for sketching polynomial kernel and proves that, for the first time, Gaussian kernel can be sketched without an exponential loss in the dimension with provable guarantees. Appendix A proves lower bounds.

2 Preliminaries

In this section we introduce notation and present useful properties of tensor product of vectors and matrices as well as properties of linear sketch matrices.

We denote the tensor product of vectors a, b by $a \otimes b$ which is formally defined as follows,

Definition 4 (Tensor product of vectors). Given $a \in \mathbb{R}^m$ and $b \in \mathbb{R}^n$ we define the *twofold tensor product* $a \otimes b$ to be

$$a \otimes b = \begin{bmatrix} a_1 b_1 & a_1 b_2 & \cdots & a_1 b_n \\ a_2 b_1 & a_2 b_2 & \cdots & a_2 b_n \\ \vdots & \vdots & & \vdots \\ a_m b_1 & a_m b_2 & \cdots & a_m b_n \end{bmatrix} \in \mathbb{R}^{m \times n}.$$

Although tensor products are multidimensional objects, it is often convenient to associate them with single-dimensional vectors. In particular, we will often associate $a \otimes b$ with the single-dimensional column vector $(a_1 b_1, a_2 b_1, \dots, a_m b_1, a_1 b_2, a_2 b_2, \dots, a_m b_2, \dots, a_m b_n)$. Given $v_1 \in \mathbb{R}^{d_1}, v_2 \in \mathbb{R}^{d_2} \dots v_k \in \mathbb{R}^{d_k}$, we define the *k-fold tensor product* $v_1 \otimes v_2 \cdots \otimes v_k \in \mathbb{R}^{d_1 d_2 \cdots d_k}$. For shorthand, we use the notation $v^{\otimes k}$ to denote $\underbrace{v \otimes v \cdots \otimes v}_{k \text{ terms}}$, the *k-fold self-tensoring* of v .

Tensor product can be naturally extended to matrices which is formally defined as follows,

Definition 5. Given $A_1 \in \mathbb{R}^{m_1 \times n_1}, A_2 \in \mathbb{R}^{m_2 \times n_2}, \dots, A_k \in \mathbb{R}^{m_k \times n_k}$, we define $A_1 \times A_2 \times \cdots \times A_k$ to be the matrix in $\mathbb{R}^{m_1 m_2 \cdots m_k \times n_1 n_2 \cdots n_k}$ whose element at row (i_1, \dots, i_k) and column (j_1, \dots, j_k) is $A_1(i_1, j_1) \cdots A_k(i_k, j_k)$. As a consequence the following holds for any $v_1 \in \mathbb{R}^{n_1}, v_2 \in \mathbb{R}^{n_2}, \dots, v_k \in \mathbb{R}^{n_k}$: $(A_1 \times A_2 \times \cdots \times A_k)(v_1 \otimes v_2 \otimes \cdots \otimes v_k) = (A_1 v_1) \otimes (A_2 v_2) \otimes \cdots \otimes (A_k v_k)$.

The tensor product has the useful *mixed product property*, given in the following Claim,

Claim 6. For every matrices A, B, C, D with appropriate sizes, the following holds,

$$(A \cdot B) \times (C \cdot D) = (A \times C) \cdot (B \times D).$$

We also define the column wise tensoring of matrices as follows,

Definition 6. Given $A_1 \in \mathbb{R}^{m_1 \times n}, A_2 \in \mathbb{R}^{m_2 \times n}, \dots, A_k \in \mathbb{R}^{m_k \times n}$, we define $A_1 \otimes A_2 \otimes \dots \otimes A_k$ to be the matrix in $\mathbb{R}^{m_1 m_2 \dots m_k \times n}$ whose j^{th} column is $A_1^j \otimes A_2^j \otimes \dots \otimes A_k^j$ for every $j \in [n]$, where A_l^j is the j^{th} column of A_l for every $l \in [k]$.

Similarly the row wise tensoring of matrices are introduced in the following Definition,

Definition 7. Given $A^1 \in \mathbb{R}^{m \times n_1}, A^2 \in \mathbb{R}^{m \times n_2}, \dots, A^k \in \mathbb{R}^{m \times n_k}$, we define $A^1 \bullet A^2 \bullet \dots \bullet A^k$ to be the matrix in $\mathbb{R}^{m \times n_1 n_2 \dots n_k}$ whose j^{th} row is $(A_1^j \otimes A_2^j \otimes \dots \otimes A_k^j)^\top$ for every $j \in [m]$, where A_l^j is the j^{th} row of A^l as a column vector for every $l \in [k]$.

Definition 8. Another related operation is the *direct sum* for vectors: $x \oplus y = \begin{bmatrix} x \\ y \end{bmatrix}$ and for matrices: $A \oplus B = \begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix}$. When the sizes match up, we have $(A \oplus B)(x \oplus y) = Ax + By$. Also notice that if I_k is the $k \times k$ identity matrix, then $I_k \otimes A = \underbrace{A \oplus \dots \oplus A}_{k \text{ times}}$.

3 Construction of the Sketch

In this section, we present the basic construction for our new sketch. Suppose we are given $v_1, v_2, \dots, v_q \in \mathbb{R}^m$. Our main task is to map the tensor product $v_1 \otimes v_2 \otimes \dots \otimes v_q$ to a vector of size m using a linear sketch.

Our sketch construction is recursive in nature. To illustrate the general idea, let us first consider the case in which $q \geq 2$ is a power of two. Our sketch involves first sketching each pair $(v_1 \otimes v_2), (v_3 \otimes v_4), \dots, (v_{q-1} \otimes v_q) \in \mathbb{R}^{m^2}$ independently using independent instances of some linear base sketch (e.g., degree two TensorSketch, Sub-sampled Randomized Hadamard Transform (SRHT), CountSketch, OSNAP). The number of vectors after this step is half of the number of vectors that we began with. The natural idea is to recursively apply the same procedure on the sketched tensors with half as many instances of the base sketch in each successive step.

More precisely, we first choose a (randomized) base sketch $S_{\text{base}} : \mathbb{R}^{m^2} \rightarrow \mathbb{R}^m$ that sketches twofold tensor products of vectors in \mathbb{R}^m (we will describe how to choose the base sketch later). Then, for any power of two $q \geq 2$, we define $Q^q : \mathbb{R}^{m^q} \rightarrow \mathbb{R}^m$ on $v_1 \otimes v_2 \otimes \dots \otimes v_q$ recursively as follows:

$$Q^q(v_1 \otimes v_2 \otimes \dots \otimes v_q) = Q^{q/2} \left(S_1^q(v_1 \otimes v_2) \otimes S_2^q(v_3 \otimes v_4) \otimes \dots \otimes S_{q/2}^q(v_{q-1} \otimes v_q) \right),$$

where $S_1^q, S_2^q, \dots, S_{q/2}^q : \mathbb{R}^{m^2} \rightarrow \mathbb{R}^m$ are independent instances of S_{base} and $Q^1 : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is simply the identity map on \mathbb{R}^m .

The above construction of Q^q has been defined in terms of its action on q -fold tensor products of vectors in \mathbb{R}^m , but it extends naturally to a linear mapping from \mathbb{R}^{m^q} to \mathbb{R}^m . The formal definition of Π^q is presented below.

Definition 9 (Sketch Q^q). Let $m \geq 2$ be a positive integer and let $S_{\text{base}} : \mathbb{R}^{m^2} \rightarrow \mathbb{R}^m$ be a linear map that specifies some base sketch. Then, for any integer power of two $q \geq 2$, we define $Q^q : \mathbb{R}^{m^q} \rightarrow \mathbb{R}^m$ to be the linear map specified as follows:

$$Q^q \equiv S^2 \cdot S^4 \dots S^{q/2} \cdot S^q,$$

where for each $l \in \{2^1, 2^2, \dots, q/2, q\}$, S^l is a matrix in $\mathbb{R}^{m^{l/2} \times m^l}$ defined as

$$S^l \equiv S_1^l \times S_2^l \times \dots \times S_{l/2}^l, \tag{2}$$

where the matrices $S_1^l, \dots, S_{l/2}^l \in \mathbb{R}^{m \times m^2}$ are drawn independently from a base distribution S_{base} .

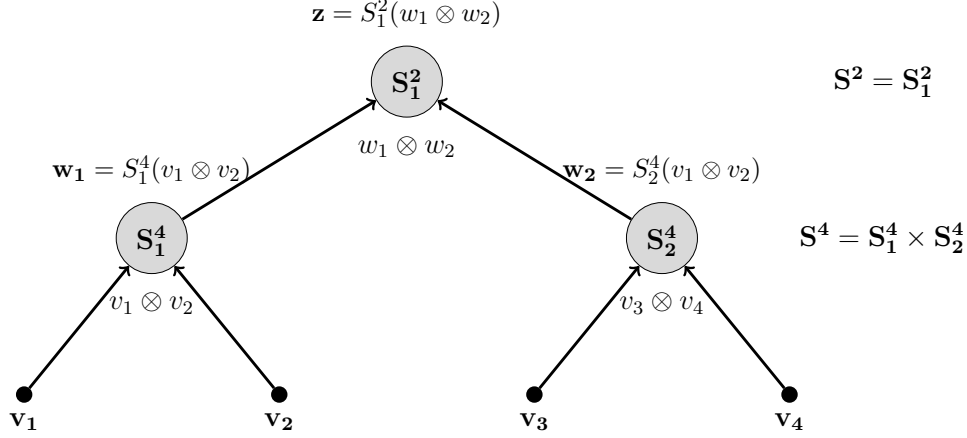


Figure 2: Visual illustration of the recursive construction of Q^q for degree $q = 4$. The input tensor is $v_1 \otimes v_2 \otimes v_3 \otimes v_4$ and the output is $z = Q^4(v_1 \otimes v_2 \otimes v_3 \otimes v_4)$. The intermediate nodes sketch the tensors $w_1 = S_1^4(v_1 \otimes v_2)$ and $w_2 = S_2^4(v_3 \otimes v_4)$.

This sketch construction can be best visualized using a balanced binary tree with q leaves. Figure 2 illustrates the construction of degree 4, Q^4 .

For every integer q which is a power of two, by definition of S^q in (2) of Definition 9, $S^q = S_1^q \times \dots \times S_{q/2}^q$. Hence, by claim 6 we can write,

$$S^q = S_1^q \times \dots \times S_{q/2}^q = \left(S_1^q \times \dots \times S_{q/2-1}^q \times I_m \right) \cdot \left(I_{m^{q-2}} \times S_{q/2}^q \right).$$

By multiple applications of Claim 6 we have the following claim,

Claim 7. For every power of two integer q and any positive integer m , if S^q is defined as in (2) of Definition 9, then

$$S^q = M_{q/2} M_{q/2-1} \dots M_1,$$

where $M_j = I_{m^{q-2j}} \times S_{q/2-j+1}^q \times I_{m^{j-1}}$ for every $j \in [q/2]$.

Embedding \mathbb{R}^{d^q} : So far we have constructed a sketch Q^q for sketching tensor product of vectors in \mathbb{R}^m . However, in general the data points can be in a space \mathbb{R}^d of arbitrary dimension. A natural idea is to reduce the dimension of the vectors by a mapping from \mathbb{R}^d to \mathbb{R}^m and then apply Q^q on the tensor product of reduced data points. The dimensionality reduction defines a linear mapping from \mathbb{R}^{d^q} to \mathbb{R}^{m^d} which can be represented by a matrix. We denote the dimensionality reduction matrix by $T^q \in \mathbb{R}^{m^d \times d^q}$ formally defined as follows.

Definition 10. Let m, d be positive integers and let $T_{\text{base}} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be a linear map that specifies some base sketch. Then for any integer power of two q we define T^q to be the linear map specified as follows,

$$T^q = T_1 \times T_2 \times \dots \times T_q,$$

where the matrices T_1, \dots, T_q are drawn independently from T_{base} .

Discussion: Similar to Claim 7, the transform T^q can be expressed as the following product of q matrices,

$$T^q = M_q M_{q-1} \dots M_1,$$

where $M_j = I_{d^{q-j}} \times T_{q-j+1} \times I_{d^{j-1}}$ for every $j \in [q]$.

Now we define the final sketch $\Pi^q : \mathbb{R}^{d^q} \rightarrow \mathbb{R}^m$ for arbitrary d as the composition of $Q^q \cdot T^q$. Moreover, to extend the definition to arbitrary q which is not necessarily a power of two we tensor the input vector with a standard basis vector a number of times to make the input size compatible with the sketch matrices. The sketch Π^q is formally defined below,

Definition 11 (Sketch Π^p). Let m, d be positive integers and let $S_{\text{base}} : \mathbb{R}^{m^2} \rightarrow \mathbb{R}^m$ and $T_{\text{base}} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ be linear maps that specify some base sketches. Then, for any integer $p \geq 2$ we define $\Pi^p : \mathbb{R}^{d^p} \rightarrow \mathbb{R}^m$ to be the linear map specified as follows:

1. If p is a power of two, then Π^p is defined as

$$\Pi^p = Q^p \cdot T^p,$$

where $Q^p \in \mathbb{R}^{m \times m^p}$ and $T^p \in \mathbb{R}^{m^p \times d^p}$ are sketches as in Definitions 9 and 10 respectively.

2. If p is not a power of two, then let $q = 2^{\lceil \log_2 p \rceil}$ be the smallest power of two integer that is greater than p and we define Π^p as

$$\Pi^p(v) = \Pi^q \left(v \otimes e_1^{\otimes (q-p)} \right),$$

for every $v \in \mathbb{R}^{d^p}$, where $e_1 \in \mathbb{R}^d$ is the standard basis column vector with a 1 in the first coordinate and zeros elsewhere, and Π^q is defined as in the first part of this definition.

Algorithm 1 sketches $x^{\otimes p}$ for any integer p and any input vector $x \in \mathbb{R}^d$ using the sketch Π^p as in Definition 11, i.e., computes $\Pi^p(x^{\otimes p})$.

Algorithm 1 SKETCH FOR THE TENSOR $x^{\otimes p}$

input: vector $x \in \mathbb{R}^d$, dimension d , degree p , number of buckets m , base sketches $S_{\text{base}} \in \mathbb{R}^{m \times m^2}$ and $T_{\text{base}} \in \mathbb{R}^{m \times d}$

output: sketched vector $z \in \mathbb{R}^m$

- 1: Let $q = 2^{\lceil \log_2 p \rceil}$
 - 2: Let T_1, \dots, T_q be independent instances of the base sketch $T_{\text{base}} : \mathbb{R}^d \rightarrow \mathbb{R}^m$
 - 3: For every $j \in \{1, 2, \dots, p\}$, let $Y_j^0 = T_j \cdot x$
 - 4: For every $j \in \{p+1, \dots, q\}$, let $Y_j^0 = T_j \cdot e_1$, where e_1 is the standard basis vector in \mathbb{R}^d with value 1 in the first coordinate and zero elsewhere
 - 5: **for** $l = 1$ to $\log_2 q$ **do**
 - 6: Let $S_1^{q/2^{l-1}}, \dots, S_{q/2^l}^{q/2^{l-1}}$ be independent instances of the base sketch $S_{\text{base}} : \mathbb{R}^{m^2} \rightarrow \mathbb{R}^m$
 - 7: For every $j \in \{1, \dots, q/2^l\}$ let $Y_j^l = S_j^{q/2^{l-1}} \left(Y_{2j-1}^{l-1} \otimes Y_{2j}^{l-1} \right)$
 - 8: **end for**
 - 9: **return** $z = Y_1^{\log_2 q}$
-

We show the correctness of Algorithm 1 in the next lemma.

Lemma 8. For any positive integers d, m , and p , any distribution on matrices $S_{\text{base}} : \mathbb{R}^{m^2} \rightarrow \mathbb{R}^m$ and $T_{\text{base}} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ which specify some base sketches, any vector $x \in \mathbb{R}^d$, Algorithm 1 computes $\Pi^p(x^{\otimes p})$ as in Definition 11.

Proof. For every input vector $x \in \mathbb{R}^d$ to Algorithm 1, the vectors Y_1^0, \dots, Y_p^0 , are computed in lines 3 and 4 of algorithm as $Y_j^0 = T_j \cdot x$, for all $j \in \{1, \dots, p\}$, and, $Y_{j'}^0 = T_{j'} \cdot e_1$, for all $j \in \{q+1, \dots, q\}$. Therefore, as shown in Definition 5, the following holds,

$$Y_1^0 \otimes \cdots \otimes Y_p^0 = T_1 \times \cdots \times T_q \cdot (x^{\otimes p} \otimes e_1^{\otimes(q-p)}).$$

From the definition of sketch T^q as per Definition 10 it follows that,

$$Y_1^0 \otimes \cdots \otimes Y_q^0 = T^q \cdot (x^{\otimes p} \otimes e_1^{\otimes(q-p)}). \quad (3)$$

The algorithm computes $Y_1^l, \dots, Y_{q/2^l}^l$ in line 7 as, $Y_j^l = S_j^{q/2^{l-1}} (Y_{2j-1}^{l-1} \otimes Y_{2j}^{l-1})$, for every $j \in \{1, \dots, q/2^l\}$ and every $l \in \{1, \dots, \log_2 q\}$ in a for loop. Therefore, by Claim 6,

$$Y_1^l \otimes \cdots \otimes Y_{q/2^l}^l = \left(S_1^{q/2^{l-1}} \times \cdots \times S_{q/2^l}^{q/2^{l-1}} \right) \cdot Y_1^{l-1} \otimes \cdots \otimes Y_{q/2^{l-1}}^{l-1}.$$

By the definition of the sketch $S^{q/2^{l-1}}$ in (2) of Definition 9 we have that for every $l \in \{1, \dots, \log_2 q\}$,

$$Y_1^l \otimes \cdots \otimes Y_{q/2^l}^l = S^{q/2^{l-1}} \cdot Y_1^{l-1} \otimes \cdots \otimes Y_{q/2^{l-1}}^{l-1}.$$

Therefore, by recursive application of the above identity we get that,

$$Y_1^{\log_2 p} = S^2 \cdot S^4 \dots S^{q/2} \cdot S^q \cdot Y_1^0 \otimes \cdots \otimes Y_q^0.$$

From the definition of sketch Q^q as in Definition 9 it follows that,

$$Y_1^{\log_2 q} = Q^q \cdot Y_1^0 \otimes \cdots \otimes Y_q^0.$$

Substituting $Y_1^0 \otimes \cdots \otimes Y_q^0$ from (3) in the above gives, $z = (Q^q \cdot T^q) \cdot (x^{\otimes p} \otimes e_1^{\otimes(q-p)})$, where by Definition 11 we have that, $z = \Pi^p(x^{\otimes p})$. \square

Choices of the Base Sketches S_{base} and T_{base} : We present formal definitions for various choices of the base sketches S_{base} and T_{base} that will be used for our sketch construction Π^q of Definition 11. We start by briefly recalling the **CountSketch** [CCFC02].

Definition 12 (CountSketch transform). Let $h : [d] \rightarrow [m]$ be a 3-wise independent hash function and also let $\sigma : [d] \rightarrow \{-1, +1\}$ be a 4-wise independent random sign function. Then, the CountSketch transform, $S : \mathbb{R}^d \rightarrow \mathbb{R}^m$, is defined as follows; for every $i \in [d]$ and every $r \in [m]$,

$$S_{r,i} = \sigma(i) \cdot \mathbb{1}[h(i) = r].$$

Another base sketch that we consider is the **TensorSketch** of degree two [Pag13] defined as follows.

Definition 13 (degree two TensorSketch transform). Let $h_1, h_2 : [d] \rightarrow [m]$ be 3-wise independent hash functions and also let $\sigma_1, \sigma_2 : [d] \rightarrow \{-1, +1\}$ be 4-wise independent random sign functions. Then, the degree two TensorSketch transform, $S : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^m$, is defined as follows; for every $i, j \in [d]$ and every $r \in [m]$,

$$S_{r,(i,j)} = \sigma_1(i) \cdot \sigma_2(j) \cdot \mathbb{1}[h_1(i) + h_2(j) = r \pmod{m}].$$

Remark: $S(x^{\otimes 2})$ can be computed in $O(m \log m + \text{nnz}(x))$ time using the Fast Fourier Transform.

Now let us briefly recall the SRHT [AC06].

Definition 14 (Subsampled Randomized Hadamard Transform (SRHT)). Let D be a $d \times d$ diagonal matrix with independent Rademacher random variables along the diagonal. Also, let $P \in \{0, 1\}^{m \times d}$ be a random sampling matrix in which each row contains a 1 at a uniformly distributed coordinate and zeros elsewhere, and let H be a $d \times d$ Hadamard matrix. Then, the SRHT, $S \in \mathbb{R}^{m \times d}$, is $S = \frac{1}{\sqrt{m}}PHD$.

We now define a variant of the SRHT which is very efficient for sketching $x^{\otimes 2}$ which we call the *TensorSRHT*.

Definition 15 (Tensor Subsampled Randomized Hadamard Transform (TensorSRHT)). Let D_1 and D_2 be two independent $d \times d$ diagonal matrices, each with diagonal entries given by independent Rademacher variables. Also let $P \in \{0, 1\}^{m \times d^2}$ be a random sampling matrix in which each row contains exactly one uniformly distributed nonzero element which has value one, and let H be a $d \times d$ Hadamard matrix. Then, the TensorSRHT is defined to be $S : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^m$ given by $S = \frac{1}{\sqrt{m}}P \cdot (HD_1 \times HD_2)$.

Remark: $S(x^{\otimes 2})$ can be computed in time $O(d \log d + m)$ using the FFT algorithm.

Another sketch which is particularly efficient for sketching sparse vectors with high probability is the OSNAP transform [NN13], defined as follows.

Definition 16 (OSNAP transform). For every sparsity parameter s , target dimension m , and positive integer d , the OSNAP transform with sparsity parameter s is defined as,

$$S_{r,j} = \sqrt{\frac{1}{s}} \cdot \delta_{r,j} \cdot \sigma_{r,j},$$

for all $r \in [m]$ and all $j \in [d]$, where $\sigma_{r,j} \in \{-1, +1\}$ are independent and uniform Rademacher random variables and $\delta_{r,j}$ are Bernoulli random variables satisfying,

1. For every $i \in [d]$, $\sum_{r \in [m]} \delta_{r,i} = s$. That is, each column of S has exactly s non-zero entries.
2. For all $r \in [m]$ and all $i \in [d]$, $\mathbb{E}[\delta_{r,i}] = s/m$.
3. The $\delta_{r,i}$'s are negatively correlated: $\forall T \subset [m] \times [d]$, $\mathbb{E}\left[\prod_{(r,i) \in T} \delta_{r,i}\right] \leq \prod_{(r,i) \in T} \mathbb{E}[\delta_{r,i}] = \left(\frac{s}{m}\right)^{|T|}$.

4 Linear Dependence on the Tensoring Degree p

There are various desirable properties that we would like a linear sketch to satisfy. One such property which is central to our main results is the *JL Moment Property*. In this section we prove Theorem 1 and Theorem 2 by propagating the *JL Moment Property* through our recursive construction from Section 3. The *JL Moment Property* captures a bound on the moments of the difference between the Euclidean norm of a vector and its Euclidean norm after applying the sketch on it. The JL Moment Property proves to be a powerful property for a sketch and we will show that it implies the Oblivious Subspace Embedding as well as the Approximate Matrix Product property for linear sketches.

In section 4.1 we choose S_{base} and T_{base} to be *TensorSketch* and *CountSketch* respectively. Then we propagate the second JL Moment through the sketch construction Π^p and thereby prove Theorem 1. In section 4.2 we propagate the higher JL Moments through our recursive construction Π^p as per Definition 11 with *TensorSRHT* at the internal nodes (S_{base}) and *OSNAP* at the leaves (T_{base}), thereby proving Theorem 2.

To make the notation less heavy we will use $\|X\|_{L^t}$ for the t^{th} moment of a random variable X . This is formally defined below.

Definition 17. For every integer $t \geq 1$ and any random variable $X \in \mathbb{R}$, we write

$$\|X\|_{L^t} = \left(E \left[|X|^t \right] \right)^{1/t}.$$

Note that $\|X + Y\|_{L^t} \leq \|X\|_{L^t} + \|Y\|_{L^t}$ for any random variables X, Y by the Minkowski Inequality.

We now formally define the JL Moment Property of sketches.

Definition 18 (JL Moment Property). For every positive integer t and every $\delta, \varepsilon \geq 0$, we say a distribution over random matrices $S \in \mathbb{R}^{m \times d}$ has the (ε, δ, t) -JL-moment property, when

$$\left\| \|Sx\|_2^2 - 1 \right\|_{L^t} \leq \varepsilon \delta^{1/t} \quad \text{and} \quad \mathbb{E} \left[\|Sx\|_2^2 \right] = 1$$

for all $x \in \mathbb{R}^d$ such that $\|x\| = 1$.

The JL Moment Property directly implies the following moment bound for the inner product of two vectors:

Lemma 9 (Two vector JL Moment Property). *For any $x, y \in \mathbb{R}^d$, if S has the (ε, δ, t) -JL Moment Property, then*

$$\left\| (Sx)^\top (Sy) - x^\top y \right\|_{L^t} \leq \varepsilon \delta^{1/t} \|x\|_2 \|y\|_2. \quad (4)$$

Proof. We can assume by linearity of the norms that $\|x\|_2 = \|y\|_2 = 1$. We then use that $\|x - y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 - 2x^\top y$ and $\|x + y\|_2^2 = \|x\|_2^2 + \|y\|_2^2 + 2x^\top y$ such that $x^\top y = (\|x + y\|_2^2 - \|x - y\|_2^2)/4$. Plugging this into the left hand side of (4) gives

$$\begin{aligned} \left\| (Sx)^\top (Sy) - x^\top y \right\|_{L^t} &= \left\| \|Sx + Sy\|_2^2 - \|x + y\|_2^2 - \|Sx - Sy\|_2^2 + \|x - y\|_2^2 \right\|_{L^t} / 4 \\ &\leq \left(\left\| \|S(x + y)\|_2^2 - \|x + y\|_2^2 \right\|_{L^t} + \left\| \|S(x - y)\|_2^2 - \|x - y\|_2^2 \right\|_{L^t} \right) / 4 \\ &\leq \varepsilon \delta^{1/t} (\|x + y\|_2^2 + \|x - y\|_2^2) / 4 \quad (\text{JL moment property}) \\ &= \varepsilon \delta^{1/t} (\|x\|_2^2 + \|y\|_2^2) / 2 \\ &= \varepsilon \delta^{1/t}. \end{aligned}$$

□

We will also need the Strong JL Moment Property, which is a sub-Gaussian bound on the difference between the Euclidean norm of a vector and its Euclidean norm after applying the sketch on it.

Definition 19 (Strong JL Moment Property). For every $\varepsilon, \delta > 0$ we say a distribution over random matrices $M \in \mathbb{R}^{m \times d}$ has the Strong (ε, δ) -JL Moment Property when

$$\left\| \|Mx\|_2^2 - 1 \right\|_{L^t} \leq \frac{\varepsilon}{e} \sqrt{\frac{t}{\log(1/\delta)}} \quad \text{and} \quad \mathbb{E} \left[\|Mx\|_2^2 \right] = 1,$$

for all $x \in \mathbb{R}^d$, $\|x\|_2 = 1$ and every integer $t \leq \log(1/\delta)$.

Remark 1. It should be noted that if a matrix $M \in \mathbb{R}^{m \times d}$ has the Strong (ε, δ) -JL Moment Property then it has the $(\varepsilon, \delta, \log(1/\delta))$ -JL Moment Property, since

$$\left\| \|Mx\|_2^2 - 1 \right\|_{L^{\log(1/\delta)}} \leq \frac{\varepsilon}{e} \sqrt{\frac{\log(1/\delta)}{\log(1/\delta)}} = \frac{\varepsilon}{e} = \varepsilon \delta^{1/\log(\frac{1}{\delta})}.$$

The following two lemmas together show that if we want to prove that Π^p is an Oblivious Subspace Embedding and that Π^p has the Approximate Matrix Multiplication Property, then it suffices to prove that Π^q has the JL Moment Property, for q which is the smallest power of two integer such that $q \geq p$, as in Definition 11. This reduction will be the main component of the proofs of Theorem 1 and Theorem 2.

Lemma 10. *For every positive integers n, p, d , every $\varepsilon, \delta \in [0, 1]$, and every $\mu \geq 0$. Let $q = 2^{\lceil \log_2(p) \rceil}$ and let $\Pi^p \in \mathbb{R}^{m \times d^p}$ and $\Pi^q \in \mathbb{R}^{m \times d^q}$ be defined as in Definition 11, for some base sketches $S_{\text{base}} \in \mathbb{R}^{m \times m^2}$ and $T_{\text{base}} \in \mathbb{R}^{d \times d}$.*

If Π^q is an $(\varepsilon, \delta, \mu, d^q, n)$ -Oblivious Subspace Embedding then Π^p is an $(\varepsilon, \delta, \mu, d^p, n)$ -Oblivious Subspace Embedding. Also if Π^q has the (ε, δ) -Approximate Matrix Multiplication Property then Π^p has the (ε, δ) -Approximate Matrix Multiplication Property.

Proof. We will prove a correspondence between Π^p and Π^q . Let $E_1 \in \mathbb{R}^{d \times n}$ be a matrix whose first row is equal to one and is zero everywhere else. By Definition 11 we have that for any matrix $A \in \mathbb{R}^{d^p \times n}$ that $\Pi^p A = \Pi^q (A \otimes E_1^{\otimes (q-p)})$. A simple calculation shows that for any matrices $A, B \in \mathbb{R}^{d^p \times n}$ then

$$(A \otimes E_1^{\otimes (q-p)})^\top (B \otimes E_1^{\otimes (q-p)}) = A^\top B \circ (E_1^{\otimes (q-p)})^\top E_1^{\otimes (q-p)} = A^\top B ,$$

where \circ denotes the Hadamard product, and the last equality follows since $(E_1^{\otimes (q-p)})^\top E_1^{\otimes (q-p)}$ is an all ones matrix. This implies that $\|A \otimes E_1^{\otimes (q-p)}\|_F = \|A\|_F$ and $s_\lambda((A \otimes E_1^{\otimes (q-p)})^\top A \otimes E_1^{\otimes (q-p)}) = s_\lambda(A^\top A)$.

Now assume that Π^q is an $(\varepsilon, \delta, \mu, n)$ -Oblivious Subspace Embedding, and let $A \in \mathbb{R}^{d^p \times n}$ and $\lambda \geq 0$ be such that $s_\lambda(A) \leq \mu$. Define $A' = A \otimes E_1^{\otimes (q-p)}$, then

$$\begin{aligned} & \Pr \left[(1 - \varepsilon)(A^\top A + \lambda I_n) \preceq (\Pi^p A)^\top \Pi^p A + \lambda I_n \preceq (1 + \varepsilon)(A^\top A + \lambda I_n) \right] \\ &= \Pr \left[(1 - \varepsilon)(A'^\top A' + \lambda I_n) \preceq (\Pi^q A')^\top \Pi^q A' + \lambda I_n \preceq (1 + \varepsilon)(A'^\top A' + \lambda I_n) \right] \\ &\geq 1 - \delta , \end{aligned}$$

where we have used that $s_\lambda(A'^\top A') = s_\lambda(A^\top A) \leq \mu$. This shows that Π^p is an $(\varepsilon, \delta, \mu, n)$ -Oblivious Subspace Embedding.

Assume that Π^q has (ε, δ) -Approximate Matrix Multiplication Property, and let $C, D \in \mathbb{R}^{d^p \times n}$. Define $C' = C \otimes E_1^{\otimes (q-p)}$ and $D' = D \otimes E_1^{\otimes (q-p)}$, then

$$\begin{aligned} & \Pr \left[\|(\Pi^p C)^\top \Pi^p D - C^\top D\|_F \geq \varepsilon \|C\|_F \|D\|_F \right] = \Pr \left[\|(\Pi^q C')^\top \Pi^q D' - C'^\top D'\|_F \geq \varepsilon \|C'\|_F \|D'\|_F \right] \\ &\leq \delta , \end{aligned}$$

where we have used that $\|C'\|_F = \|C\|_F$, $\|D'\|_F = \|D\|_F$, and $C'^\top D' = C^\top D$. This show that Π^p has (ε, δ) -Approximate Matrix Multiplication Property. \square

Lemma 11. *For any $\varepsilon, \delta \in [0, 1]$, $t \geq 1$, if $M \in \mathbb{R}^{m \times d}$ is a random matrix with (ε, δ, t) -JL Moment Property then M has the (ε, δ) -Approximate Matrix Multiplication Property.*

Furthermore, for any $\mu > 0$, if $M \in \mathbb{R}^{m \times d}$ is a random matrix with $(\varepsilon/\mu, \delta, t)$ -JL Moment Property then for every positive integer $n \in \mathbb{Z}$, M is a $(\varepsilon, \delta, \mu, d, n)$ -OSE.

Proof.

Approximate Matrix Multiplication Let $C, D \in \mathbb{R}^{d \times n}$. We will prove that

$$\left\| \|(MC)^\top MD - C^\top D\|_F \right\|_{L^t} \leq \varepsilon \delta^{1/t} \|C\|_F \|D\|_F. \quad (5)$$

Then Markov's inequality will give us the result. Using the triangle inequality together with Lemma 9 we get that:

$$\begin{aligned} \left\| \|(MC)^\top MD - C^\top D\|_F \right\|_{L^t} &= \left\| \|(MC)^\top MD - C^\top D\|_F^2 \right\|_{L^{t/2}}^{1/2} \\ &= \left\| \sum_{i,j \in [n]} \left((MC_i)^\top MD_j - C_i^\top D_j \right)^2 \right\|_{L^{t/2}}^{1/2} \\ &\leq \sqrt{\sum_{i,j \in [n]} \|(MC_i)^\top MD_j - C_i^\top D_j\|_{L^t}^2} \\ &\leq \sqrt{\sum_{i,j \in [n]} \varepsilon^2 \delta^{2/t} \|C_i\|_2^2 \|D_j\|_2^2} \\ &= \varepsilon \delta^{1/t} \|C\|_F \|D\|_F. \end{aligned}$$

Using Markov's inequality we now get that

$$\Pr \left[\|(MC)^\top MD - C^\top D\|_F \geq \varepsilon \|C\|_F \|D\|_F \right] \leq \frac{\left\| \|(MC)^\top MD - C^\top D\|_F \right\|_{L^t}^t}{\varepsilon^t \|C\|_F^t \|D\|_F^t} \leq \delta.$$

Oblivious Subspace Embedding. We will prove that for any $\lambda \geq 0$ and any matrix $A \in \mathbb{R}^{d \times n}$,

$$(1 - \varepsilon)(A^\top A + \lambda I_n) \preceq (MA)^\top MA + \lambda I_n \preceq (1 + \varepsilon)(A^\top A + \lambda I_n), \quad (6)$$

holds with probability at least $1 - \left(\frac{s_\lambda(A^\top A)}{\mu} \right)^t \delta$, which will imply our result.

We will first consider $\lambda > 0$. Then $A^\top A + \lambda I_n$ is positive definite. Thus, by left and right multiplying (6) by $(A^\top A + \lambda I_n)^{-1/2}$, we see that (6) is equivalent to

$$(1 - \varepsilon)I_n \preceq \left(MA(A^\top A + \lambda I_n)^{-1/2} \right)^\top MA(A^\top A + \lambda I_n)^{-1/2} + \lambda(A^\top A + \lambda I_n)^{-1} \preceq (1 + \varepsilon)I_n.$$

which, in turn, is implied by the following:

$$\left\| \left(MA(A^\top A + \lambda I_n)^{-1/2} \right)^\top MA(A^\top A + \lambda I_n)^{-1/2} + \lambda(A^\top A + \lambda I_n)^{-1} - I_n \right\|_{op} \leq \varepsilon.$$

Note that $(A^\top A + \lambda I_n)^{-1/2} A^\top A (A^\top A + \lambda I_n)^{-1/2} = I_n - \lambda(A^\top A + \lambda I_n)^{-1}$. Letting $Z = A(A^\top A + \lambda I_n)^{-1/2}$, we note that it suffices to establish,

$$\left\| (MZ)^\top MZ - Z^\top Z \right\|_{op} \leq \varepsilon.$$

Using (5) together with Markov's inequality we get that

$$\Pr \left[\left\| (MZ)^\top MZ - Z^\top Z \right\|_{op} \geq \varepsilon \right] \leq \Pr \left[\left\| (MZ)^\top MZ - Z^\top Z \right\|_F \geq \varepsilon \right] \leq \left(\frac{\|Z\|_F^2}{\mu} \right)^t \delta = \left(\frac{s_\lambda(A^\top A)}{\mu} \right)^t \delta,$$

where the last equality follows from

$$\begin{aligned}
\|Z\|_F^2 &= \mathbf{tr} \left(Z^\top Z \right) \\
&= \mathbf{tr} \left(\left(A(A^\top A + \lambda I_n)^{-1/2} \right)^\top A(A^\top A + \lambda I_n)^{-1/2} \right) \\
&= \mathbf{tr} \left(A^\top A (A^\top A + \lambda I_n)^{-1} \right) \\
&= s_\lambda(A^\top A).
\end{aligned}$$

To prove the result for $\lambda = 0$ we will use Fatou's lemma.

$$\begin{aligned}
&\Pr \left[\left((1 - \varepsilon)A^\top A \preceq (MA)^\top MA \preceq (1 + \varepsilon)A^\top A \right)^C \right] \\
&\leq \liminf_{\lambda \rightarrow 0^+} \Pr \left[\left((1 - \varepsilon)(A^\top A + \lambda I_n) \preceq (MA)^\top MA + \lambda I_n \preceq (1 + \varepsilon)(A^\top A + \lambda I_n) \right)^C \right] \\
&\leq \liminf_{\lambda \rightarrow 0^+} \frac{s_\lambda(A^\top A)}{\mu} \delta \\
&= \frac{s_0(A^\top A)}{\mu} \delta,
\end{aligned}$$

where the last equality follows from continuity of $\lambda \mapsto s_\lambda(A^\top A)$. \square

Our next important observation is that Π^q can be written as the product of $2q - 1$ independent random matrices, which all have a special structure which makes them easy to analyse.

Lemma 12. *For any integer q which is a power of two, $\Pi^q : \mathbb{R}^{m^q} \rightarrow \mathbb{R}^m$ be defined as in Definition 11 for some base sketches $S_{\text{base}} : \mathbb{R}^{m^2} \rightarrow \mathbb{R}^m$ and $T_{\text{base}} : \mathbb{R}^d \rightarrow \mathbb{R}^m$. Then there exist matrices $(M^{(i)})_{i \in [q-1]}$, $(M'^{(j)})_{j \in [q]}$ and integers $(k_i)_{i \in [q-1]}$, $(k'_i)_{i \in [q-1]}$, $(l_j)_{j \in [q]}$, $(l'_j)_{j \in [q]}$, such that,*

$$\Pi^q = M^{(q-1)} \cdot \dots \cdot M^{(1)} \cdot M'^{(q)} \cdot \dots \cdot M'^{(1)},$$

and $M^{(i)} = I_{k_i} \times S_{\text{base}}^{(i)} \times I_{k'_i}$, $M'^{(j)} = I_{l_j} \times T_{\text{base}}^{(j)} \times I_{l'_j}$, where $S_{\text{base}}^{(i)}$ and $T_{\text{base}}^{(j)}$ are independent instances of S_{base} and T_{base} , for every $i \in [q - 1]$, $j \in [q]$.

Proof. We have that $\Pi^q = Q^q T^q$ by Definition 11. By Definition 9 we have that $Q^q = S^2 S^4 \dots S^q$. Claim 7 shows that for every $l \in \{2, 4, \dots, q\}$ we can write,

$$S^l = M_{l/2}^l M_{l/2-1}^l \cdot \dots \cdot M_1^l, \quad (7)$$

where $M_j^l = I_{m^{l-2j}} \times S_{l/2-j+1}^l \times I_{m^{j-1}}$ for every $j \in [l/2]$. From the discussion in Definition 10 it follows that,

$$T^q = M'^{(q)} \cdot \dots \cdot M'^{(1)}, \quad (8)$$

where $M'^{(j)} = I_{d^{q-j}} \times T_{q-j+1} \times I_{m^{j-1}}$ for every $j \in [q]$. Therefore by combining (7) and (8) we get the result. \square

We want to show that $I_k \times M \times I_{k'}$ inherits the JL properties of M . The following simple fact does just that.

Lemma 13. Let $t \in \mathbb{N}$ and $\alpha \geq 0$. If $P \in \mathbb{R}^{m_1 \times d_1}$ and $Q \in \mathbb{R}^{m_2 \times d_2}$ are two random matrices (not necessarily independent), such that,

$$\begin{aligned} \left\| \|Px\|_2^2 - \|x\|_2^2 \right\|_{L^t} &\leq \alpha \|x\|_2^2 \quad \text{and} \quad \mathbb{E} \left[\|Px\|_2^2 \right] = \|x\|_2^2, \\ \left\| \|Qy\|_2^2 - \|y\|_2^2 \right\|_{L^t} &\leq \alpha \|y\|_2^2 \quad \text{and} \quad \mathbb{E} \left[\|Qy\|_2^2 \right] = \|y\|_2^2, \end{aligned}$$

for any vectors $x \in \mathbb{R}^{d_1}$ and $y \in \mathbb{R}^{d_2}$, then

$$\left\| \|(P \oplus Q)z\|_2^2 - \|z\|_2^2 \right\|_{L^t} \leq \alpha \|z\|_2^2 \quad \text{and} \quad \mathbb{E} \left[\|(P \oplus Q)z\|_2^2 \right] = \|z\|_2^2,$$

for any vector $z \in \mathbb{R}^{d_1+d_2}$.

Proof. Let $z \in \mathbb{R}^{d_1+d_2}$ and choose $x \in \mathbb{R}^{d_1}$ and $y \in \mathbb{R}^{d_2}$, such that, $z = x \oplus y$. Using the triangle inequality,

$$\begin{aligned} \left\| \|(P \oplus Q)z\|_2^2 - \|z\|_2^2 \right\|_{L^t} &= \left\| \|Px\|_2^2 + \|Qy\|_2^2 - \|x\|_2^2 - \|y\|_2^2 \right\|_{L^t} \\ &\leq \left\| \|Px\|_2^2 - \|x\|_2^2 \right\|_{L^t} + \left\| \|Qy\|_2^2 - \|y\|_2^2 \right\|_{L^t} \\ &\leq \alpha \|x\|_2^2 + \alpha \|y\|_2^2 \\ &= \alpha \|z\|_2^2. \end{aligned}$$

We also see that

$$\mathbb{E} \left[\|(P \oplus Q)z\|_2^2 \right] = \mathbb{E} \left[\|Px\|_2^2 \right] + \mathbb{E} \left[\|Qy\|_2^2 \right] = \|x\|_2^2 + \|y\|_2^2 = \|z\|_2^2.$$

□

An easy consequence of this lemma is that for any matrix, S , with the (ε, δ, t) -JL Moment Property, $I_k \times S$ has the (ε, δ, t) -JL Moment Property. This follows simply from $I_k \times S = \underbrace{S \oplus S \oplus \dots \oplus S}_{k \text{ times}}$.

Similarly, $S \times I_k$ has the (ε, δ, t) -JL Moment Property, since $S \times I_k$ is just a reordering of the rows of $I_k \times S$, which trivially does not affect the JL Moment Property. The same arguments show that if S has the Strong (ε, δ) -JL Moment Property then $I_k \times S$ and $S \times I_k$ has the Strong (ε, δ) -JL Moment Property. So we conclude the following

Lemma 14. If the matrix S has the (ε, δ, t) -JL Moment Property, then for any positive integers k, k' , the matrix $M = I_k \times S \times I_{k'}$ has the (ε, δ, t) -JL Moment Property.

Similarly, if the matrix S has the Strong (ε, δ) -JL Moment Property, then for any positive integers k, k' , the matrix $M = I_k \times S \times I_{k'}$ has the Strong (ε, δ) -JL Moment Property.

Now if we can prove that the product of matrices with the JL Moment Property has the JL Moment Property, then Lemma 14 and Lemma 12 would imply that Π^q has the JL Moment Property, which again implies that Π^p is an Oblivious Subspace Embedding and has the Approximate Matrix Multiplication Property, by Lemma 11 and Lemma 10. This is exactly what we will do: in Section 4.1 we prove that the product of k independent matrices with the $(\frac{\varepsilon}{\sqrt{2k}}, \delta, 2)$ -JL Moment Property results in a matrix with the $(\varepsilon, \delta, 2)$ -JL Moment Property, which will give us the proof of Theorem 1, and in Section 4.2 we prove that the product of k independent matrices with the Strong $(O(\frac{\varepsilon}{\sqrt{k}}), \delta)$ -JL Moment Property results in a matrix with the Strong (ε, δ) -JL Moment Property, which will give us the proof of Theorem 2.

4.1 Second Moment of Π^q (analysis for T_{base} : CountSketch and S_{base} : TensorSketch)

In this section we prove Theorem 1 by instantiating our recursive construction from Section 3 with **CountSketch** at the leaves and **TensorSketch** at the internal nodes of the tree. The proof proceeds by showing the second moment property – i.e., $(\varepsilon, \delta, 2)$ -JL Moment Property, for our recursive construction. We prove that our sketch Π^q satisfies the $(\varepsilon, \delta, 2)$ -JL Moment Property as per Definition 18 as long as the base sketches $S_{\text{base}}, T_{\text{base}}$ are chosen from a distribution which satisfies the second moment property. We show that this is the case for **CountSketch** and **TensorSketch**.

Lemma 14 together with Lemma 12 show that if the base sketches $S_{\text{base}}, T_{\text{base}}$ have the JL Moment Property then Π^q is the product of $2q - 1$ independent random matrices with the JL Moment Property. Therefore, understanding how matrices with the JL Moment Property compose is crucial. The following lemma shows that composing independent random matrices which have the JL Moment Property results in matrix which has the JL Moment Property with a small loss in the parameters.

Lemma 15 (Composition lemma for the second moment). *For any $\varepsilon, \delta \geq 0$ and any integer k if $M^{(1)} \in \mathbb{R}^{d_2 \times d_1}, \dots, M^{(k)} \in \mathbb{R}^{d_{k+1} \times d_k}$ are independent random matrices with the $(\frac{\varepsilon}{\sqrt{2k}}, \delta, 2)$ -JL-moment property then the product matrix $M = M^{(k)} \dots M^{(1)}$ satisfies the $(\varepsilon, \delta, 2)$ -JL-moment property.*

Proof. Let $x \in \mathbb{R}^{d_1}$ be a fixed unit norm vector. We note that for any $i \in [k]$ we have that

$$\mathbb{E} \left[\|M^{(i)} \cdot \dots \cdot M^{(1)} x\|_2^2 \mid M^{(1)}, \dots, M^{(i-1)} \right] = \|M^{(i-1)} \cdot \dots \cdot M^{(1)} x\|_2^2. \quad (9)$$

Now we will prove by induction on $i \in [k]$ that,

$$\text{Var} \left[\|M^{(i)} \cdot \dots \cdot M^{(1)} x\|_2^2 \right] \leq \left(1 + \frac{\varepsilon^2 \delta}{2k} \right)^i - 1. \quad (10)$$

For $i = 1$ the result follows from the fact that $M^{(1)}$ has the $(\varepsilon/\sqrt{2k}, \delta, 2)$ -JL moment property. Now assume that (10) is true for $i - 1$. By the law of total variance we get that

$$\begin{aligned} \text{Var} \left[\|M^{(i)} \cdot \dots \cdot M^{(1)} x\|_2^2 \right] &= \mathbb{E} \left[\text{Var} \left[\|M^{(i)} \cdot \dots \cdot M^{(1)} x\|_2^2 \mid M^{(1)}, \dots, M^{(i-1)} \right] \right] \\ &\quad + \text{Var} \left[\mathbb{E} \left[\|M^{(i)} \cdot \dots \cdot M^{(1)} x\|_2^2 \mid M^{(1)}, \dots, M^{(i-1)} \right] \right] \end{aligned} \quad (11)$$

Using (9) and the induction hypothesis we get that,

$$\begin{aligned} \text{Var} \left[\mathbb{E} \left[\|M^{(i)} \cdot \dots \cdot M^{(1)} x\|_2^2 \mid M^{(1)}, \dots, M^{(i-1)} \right] \right] &= \text{Var} \left[\|M^{(i-1)} \cdot \dots \cdot M^{(1)} x\|_2^2 \right] \\ &\leq \left(1 + \frac{\varepsilon^2 \delta}{2k} \right)^{i-1} - 1. \end{aligned} \quad (12)$$

Using that $M^{(i)}$ has the $(\varepsilon/\sqrt{2k}, \delta, 2)$ -JL moment property, (9), and the induction hypothesis we

get that,

$$\begin{aligned}
& \mathbb{E} \left[\text{Var} \left[\left\| M^{(i)} \cdot \dots \cdot M^{(1)} x \right\|_2^2 \mid M^{(1)}, \dots, M^{(i-1)} \right] \right] \\
& \leq \mathbb{E} \left[\frac{\varepsilon^2}{2k} \delta \left\| M^{(i-1)} \cdot \dots \cdot M^{(1)} x \right\|_2^4 \right] \\
& = \frac{\varepsilon^2 \delta}{2k} \left(\text{Var} \left[\left\| M^{(i-1)} \cdot \dots \cdot M^{(1)} x \right\|_2^2 \right] + \mathbb{E} \left[\left\| M^{(i-1)} \cdot \dots \cdot M^{(1)} x \right\|_2^2 \right]^2 \right) \\
& \leq \frac{\varepsilon^2 \delta}{2k} \left(\left(1 + \frac{\varepsilon^2 \delta}{2k} \right)^{i-1} - 1 + 1 \right) = \frac{\varepsilon^2 \delta}{2k} \left(1 + \frac{\varepsilon^2 \delta}{2k} \right)^{i-1}. \tag{13}
\end{aligned}$$

Plugging (12) and (13) into (11) gives,

$$\text{Var} \left[\left\| M^{(i)} \cdot \dots \cdot M^{(1)} x \right\|_2^2 \right] \leq \frac{\varepsilon^2 \delta}{2k} \left(1 + \frac{\varepsilon^2 \delta}{2k} \right)^{i-1} + \left(1 + \frac{\varepsilon^2 \delta}{2k} \right)^{i-1} - 1 = \left(1 + \frac{\varepsilon^2 \delta}{2k} \right)^i - 1.$$

Hence,

$$\text{Var} \left[\left\| Mx \right\|_2^2 \right] \leq \left(1 + \frac{\varepsilon^2 \delta}{2k} \right)^k - 1 \leq \exp(\varepsilon^2 \delta / 2) - 1 \leq \varepsilon^2 \delta,$$

which proves that M has the $(\varepsilon, \delta, 2)$ -JL moment property. \square

Equipped with the composition lemma for the second moment, we now establish the second moment property for our recursive sketch Π^q :

Corollary 16 (Second moment property for Π^q). *For any power of two integer q let $\Pi^q : \mathbb{R}^{m^q} \rightarrow \mathbb{R}^m$ be defined as in Definition 11, where both of the common distributions $S_{\text{base}} : \mathbb{R}^{m^2} \rightarrow \mathbb{R}^m$ and $T_{\text{base}} : \mathbb{R}^d \rightarrow \mathbb{R}^m$, satisfy the $(\frac{\varepsilon}{\sqrt{4q+2}}, \delta, 2)$ -JL-moment property. Then it follows that Π^q satisfies the $(\varepsilon, \delta, 2)$ -JL-moment property.*

Proof. This follows from Lemma 12 and Lemma 15. \square

Now we are ready to prove Theorem 1. Recall that $k(x, y) = \langle x, y \rangle^q$ is the polynomial kernel of degree q . One can see that $k(x, y) = \langle x^{\otimes q}, y^{\otimes q} \rangle$. Let $x_1, x_2, \dots, x_n \in \mathbb{R}^m$ be an arbitrary dataset of n points in \mathbb{R}^m . We represent the data points by matrix $X \in \mathbb{R}^{m \times n}$ whose i^{th} column is the vector x_i . Let $A \in \mathbb{R}^{m^q \times n}$ be the matrix whose i^{th} column is $x_i^{\otimes q}$ for every $i \in [n]$. For any regularization parameter $\lambda > 0$, the statistical dimension of $A^\top A$ is defined as $s_\lambda := \text{tr} \left((A^\top A)(A^\top A + \lambda I_n)^{-1} \right)$.

Theorem 1. *For every positive integers n, p, d , every $\varepsilon, s_\lambda > 0$, there exists a distribution on linear sketches $\Pi^p \in \mathbb{R}^{m \times dp}$ such that: (1) If $m = \Omega(p s_\lambda^2 \varepsilon^{-2})$, then Π^p is an $(\varepsilon, 1/10, s_\lambda, dp, n)$ -oblivious subspace embedding as in Definition 2. (2) If $m = \Omega(p \varepsilon^{-2})$, then Π^p has the $(\varepsilon, 1/10)$ -approximate matrix product property as in Definition 3.*

Moreover, for any $X \in \mathbb{R}^{d \times n}$, if $A \in \mathbb{R}^{dp \times n}$ is the matrix whose columns are obtained by the p -fold self-tensoring of each column of X then the matrix $\Pi^p A$ can be computed using Algorithm 1 in time $\tilde{O}(pnm + p \text{nnz}(X))$.

Proof. Throughout the proof, let $\delta = \frac{1}{10}$ denote the failure probability, let $q = 2^{\lceil \log_2 p \rceil}$, and let $e_1 \in \mathbb{R}^d$ be the column vector with a 1 in the first coordinate and zeros elsewhere. Let $\Pi^p \in \mathbb{R}^{m \times d^p}$ be the sketch defined in Definition 11, where the base distributions $S_{\text{base}} \in \mathbb{R}^{m \times m^2}$ and $T_{\text{base}} \in \mathbb{R}^{m \times d}$ are respectively the standard **TensorSketch** of degree two and standard **CountSketch**. It is shown in [ANW14] and [CW17] that for these choices of base sketches, S_{base} and T_{base} are both unbiased and satisfy the $(\frac{\varepsilon}{\sqrt{4q+2}}, \delta, 2)$ -JL-moment property as long as $m = \Omega(\frac{q}{\varepsilon^2 \delta})$ (see Definition 18).

Oblivious Subspace Embedding Let $m = \Omega\left(\frac{qs_\lambda^2}{\delta \varepsilon^2}\right)$ be an integer. Then S_{base} and T_{base} has the $(\frac{\varepsilon}{\sqrt{4q+2s_\lambda}}, \delta, 2)$ -JL Moment Property. Thus using Corollary 16 we conclude that Π^q has the $(\frac{\varepsilon}{s_\lambda}, \delta, 2)$ -JL Moment Property. Thus, Lemma 11 implies that Π^q is an $(\varepsilon, \delta, s_\lambda, d^q, n)$ -Oblivious Subspace Embedding, and by Lemma 10 we get that Π^p is an $(\varepsilon, \delta, s_\lambda, d^p, n)$ -Oblivious Subspace Embedding.

Approximate Matrix Multiplication. Let $m = \Omega\left(\frac{q}{\delta \varepsilon^2}\right)$. Then S_{base} and T_{base} have the $(\frac{\varepsilon}{\sqrt{4q+2}}, \delta, 2)$ -JL Moment Property. Thus, using Corollary 16 we conclude that Π^q has the $(\varepsilon, \delta, 2)$ -JL Moment Property. Thus, Lemma 11 implies that Π^q has the (ε, δ) -Approximate Matrix Multiplication Property, and by Lemma 10 we get that Π^p has the (ε, δ) -Approximate Matrix Multiplication Property.

Runtime of Algorithm 1 when the base sketch S_{base} is TensorSketch of degree two and T_{base} is CountSketch: We compute the time of running Algorithm 1 on a vector x . Computing Y_j^0 for each j in lines 3 and 4 of algorithm requires applying a **CountSketch** on either x or e_1 which takes time $O(\text{nnz}(x))$. Therefore computing all Y_j^0 's takes time $O(q \cdot \text{nnz}(x))$.

Computing each of Y_j^l 's for $l \geq 1$ in line 7 of Algorithm 1 amounts to applying a degree two **TensorSketch** of input dimension m^2 and target dimension of m on $Y_{2j-1}^{l-1} \otimes Y_{2j}^{l-1}$. This takes time $O(m \log m)$. Therefore computing Y_j^l for all $l, j \geq 1$ takes time $O(q \cdot m \log m)$. Note that $q \leq 2p$ and hence the total running time of Algorithm 1 on one vector x is $O(p \cdot m \log_2 m + p \cdot \text{nnz}(w))$. Sketching n columns of a matrix $X \in \mathbb{R}^{d \times n}$ takes time $O(p(nm \log_2 m + \text{nnz}(X)))$. \square

4.2 Higher Moments of Π^q (analysis for T_{base} : OSNAP and S_{base} : TensorSRHT)

In this section we prove Theorem 2 by instantiating our recursive construction of Section 3 with OSNAP at the leaves and TensorSRHT at the internal nodes.

The proof proceeds by showing the Strong JL Moment Property for our sketch Π^q . If a sketch satisfies the Strong JL Moment Property then it straightforwardly is an OSE and has the approximate matrix product property. This section has two goals: first is to show that SRHT, and TensorSRHT as well as OSNAP transform all satisfy the Strong JL Moment Property. The second goal of this section is to prove that our sketch construction Π^q inherits the strong JL moment property from the base sketches $S_{\text{base}}, T_{\text{base}}$.

In this section we will need Khintchine's inequality.

Lemma 17 (Khintchine's inequality [HM07]). *Let t be a positive integer, $x \in \mathbb{R}^d$, and $(\sigma_i)_{i \in [d]}$ be independent Rademacher ± 1 random variables. Then*

$$\left\| \sum_{i=1}^d \sigma_i x_i \right\|_{L^t} \leq C_t \|x\|_2,$$

where $C_t \leq \sqrt{2} \left(\frac{\Gamma((t+1)/2)}{\sqrt{\pi}} \right)^{1/t} \leq \sqrt{t}$ for all $t \geq 1$.

One may replace (σ_i) with an arbitrary independent sequence of random variables (ς_i) with $\mathbb{E}[\varsigma_i] = 0$ and $\|\varsigma_i\|_{L^r} \leq \sqrt{r}$ for any $1 \leq r \leq t$, and the lemma still holds up to a universal constant factor on the r.h.s.

First we note that the OSNAP transform satisfies the strong JL moment property.

Lemma 18. *There exists a universal constant L , such that, the following holds. Let $M \in \mathbb{R}^{m \times d}$ be a OSNAP transform with sparsity parameter s . Let $x \in \mathbb{R}^d$ be any vector with $\|x\|_2 = 1$ and $t \geq 1$, then*

$$\left\| \|Mx\|_2^2 - 1 \right\|_{L^t} \leq L \left(\sqrt{\frac{t}{m}} + \frac{t}{s} \right). \quad (14)$$

Setting $m = \Omega(\varepsilon^{-2} \log(1/\delta))$ and $s = \Omega(\varepsilon^{-1} \log(1/\delta))$ then M has the Strong (ε, δ) -JL Moment Property (Definition 19).

Proof. The proof of (14) follows from analysis in [CJN18]. They only prove it for $t = \log(1/\delta)$ but their proof is easily extended to the general case.

Now if we set $m = 4L^2 e^2 \cdot \varepsilon^{-2} \log(1/\delta)$ and $s = 2Le \cdot \varepsilon^{-1} \log(1/\delta)$ then we get that

$$\left\| \|Mx\|_2^2 - 1 \right\|_{L^t} \leq L \sqrt{\frac{t}{4L^2 e^2 \cdot \varepsilon^{-2} \log(1/\delta)}} + L \frac{t}{2Le \cdot \varepsilon^{-1} \log(1/\delta)} \leq \frac{\varepsilon}{e} \sqrt{\frac{t}{\log(1/\delta)}},$$

for every $1 \leq t \leq \log(1/\delta)$, which proves the result. \square

We continue by proving that SRHT and TensorSRHT sketches satisfy the strong JL moment property. We will do this by proving that a more general class of matrices satisfies the strong JL moment property. More precisely, let $k \in \mathbb{Z}_{>0}$ be a positive integer and $(D^{(i)})_{i \in [k]} \in \prod_{i \in [k]} \mathbb{R}^{d_i \times d_i}$ be independent matrices, each with diagonal entries given by independent Rademacher variables. Let $d = \prod_{i \in [k]} d_i$, and $P \in \{0, 1\}^{m \times d}$ be a random sampling matrix in which each row contains exactly one uniformly distributed nonzero element which has value one. Then we will prove that the matrix $M = \frac{1}{\sqrt{m}} PH(D_1 \times \dots \times D_k)$ satisfies the strong JL moment property, where H is a $d \times d$ Hadamard matrix. If $k = 1$ then M is just a SRHT, and if $k = 2$ then M is a TensorSRHT.

In order to prove this result we need a couple of lemmas. The first lemma can be seen as a version of Khintchine's inequality for higher order chaos.

Lemma 19. *Let $t \geq 1$, $k \in \mathbb{Z}_{>0}$, and $(\sigma^{(i)})_{i \in [k]} \in \prod_{i \in [k]} \mathbb{R}^{d_i}$ be independent vectors each satisfying the Khintchine inequality $\left\| \langle \sigma^{(i)}, x \rangle \right\|_{L^t} \leq C_t \|x\|_2$ for $t \geq 1$ and any vector $x \in \mathbb{R}^{d_i}$. Let $(a_{i_1, \dots, i_k})_{i_1 \in [d_1], \dots, i_k \in [d_k]}$ be a tensor in $\mathbb{R}^{d_1 \times \dots \times d_k}$, then*

$$\left\| \sum_{i_1 \in [d_1], \dots, i_k \in [d_k]} \left(\prod_{j \in [k]} \sigma_{i_j}^{(j)} \right) a_{i_1, \dots, i_k} \right\|_{L^t} \leq C_t^k \left(\sum_{i_1 \in [d_1], \dots, i_k \in [d_k]} a_{i_1, \dots, i_k}^2 \right)^{1/2},$$

for $t \geq 1$. Or, considering $a \in \mathbb{R}^{d_1 \cdots d_k}$ a vector, then simply $\left\| \langle \sigma^{(1)} \otimes \dots \otimes \sigma^{(k)}, a \rangle \right\|_{L^t} \leq C_t^k \|a\|_2$, for $t \geq 1$.

This is related to Latała's estimate for Gaussian chaoses [Lat06], but more simple in the case where a is not assumed to have special structure. Note that this implies the classical bound on the fourth moment of products of 4-wise independent hash functions [BCL⁺10, IM08, PT12], since $C_4 = 3^{1/4}$ for Rademachers we have $\mathbb{E}[\langle \sigma^{(1)} \otimes \cdots \otimes \sigma^{(k)}, x \rangle^4] \leq 3^k \|x\|_2^4$ for four-wise independent $(\sigma^{(i)})_{i \in [k]}$.

Proof. The proof will be by induction on k . For $k = 1$ then the result is by assumption. So assume that the result is true for every value up to $k - 1$. Let $B_{i_1, \dots, i_{k-1}} = \sum_{i_k \in [d_k]} \sigma_{i_k}^{(k)} a_{i_1, \dots, i_k}$. We then pull it out of the left hand term in the theorem:

$$\begin{aligned} \left\| \sum_{i_1 \in [d_1], \dots, i_c \in [d_c]} \left(\prod_{j \in [k]} \sigma_{i_j}^{(j)} \right) a_{i_1, \dots, i_k} \right\|_{L^t} &= \left\| \sum_{i_1 \in [d_1], \dots, i_{k-1} \in [d_{k-1}]} \left(\prod_{j \in [k-1]} \sigma_{i_j}^{(j)} \right) B_{i_1, \dots, i_{k-1}} \right\|_{L^t} \\ &\leq C_t^{k-1} \left\| \left(\sum_{i_1 \in [d_1], \dots, i_{k-1} \in [d_{k-1}]} B_{i_1, \dots, i_{k-1}}^2 \right)^{1/2} \right\|_{L^t} \end{aligned} \quad (15)$$

$$\begin{aligned} &= C_t^{k-1} \left\| \sum_{i_1 \in [d_1], \dots, i_{k-1} \in [d_{k-1}]} B_{i_1, \dots, i_{k-1}}^2 \right\|_{L^{t/2}}^{1/2} \\ &\leq C_t^{k-1} \left(\sum_{i_1 \in [d_1], \dots, i_{k-1} \in [d_{k-1}]} \left\| B_{i_1, \dots, i_{k-1}}^2 \right\|_{L^{t/2}} \right)^{1/2} \quad (16) \\ &= C_t^{k-1} \left(\sum_{i_1 \in [d_1], \dots, i_{k-1} \in [d_{k-1}]} \left\| B_{i_1, \dots, i_{k-1}} \right\|_{L^t}^2 \right)^{1/2}. \end{aligned}$$

Here (15) is the inductive hypothesis and (16) is the triangle inequality. It remains to bound $\left\| B_{i_1, \dots, i_{k-1}} \right\|_{L^t}^2 \leq C_t^2 \sum_{i_c \in [d_c]} a_{i_1, \dots, i_c}^2$ by Khintchine's inequality, which finishes the induction step and hence the proof. \square

The next lemma we will be using is a type of Rosenthal inequality, but which mixes large and small moments in a careful way. It bears similarity to the one sided bound in [BLM13] (Theorem 15.10) derived from the Efron Stein inequality, and the literature has many similar bounds, but we still include a proof here based on first principles.

Lemma 20. *There exists a universal constant L , such that, for $t \geq 1$ if X_1, \dots, X_k are independent non-negative random variables with t -moment, then*

$$\left\| \sum_{i \in [k]} (X_i - \mathbb{E}[X_i]) \right\|_{L^t} \leq L \left(\sqrt{t} \left\| \max_{i \in [k]} X_i \right\|_{L^t}^{1/2} \sqrt{\sum_{i \in [k]} \mathbb{E}[X_i]} + t \left\| \max_{i \in [k]} X_i \right\|_{L^t} \right).$$

Proof. Throughout these calculations L_1 , L_2 and L_3 will be universal constants.

$$\begin{aligned}
\left\| \sum_{i \in [k]} (X_i - \mathbb{E}[X_i]) \right\|_{L^t} &\leq L_1 \left\| \sum_{i \in [k]} \sigma_i X_i \right\|_{L^t} \quad (\text{Symmetrization}) \\
&\leq L_2 \sqrt{t} \left\| \sum_{i \in [k]} X_i^2 \right\|_{L^{t/2}}^{1/2} \quad (\text{Khintchine's inequality}) \\
&\leq L_2 \sqrt{t} \left\| \max_{i \in [k]} X_i \cdot \sum_{i \in [k]} X_i \right\|_{L^{t/2}}^{1/2} \quad (\text{Non-negativity}) \\
&\leq L_2 \sqrt{t} \left\| \max_{i \in [k]} X_i \right\|_{L^t}^{1/2} \cdot \left\| \sum_{i \in [k]} X_i \right\|_{L^t}^{1/2} \quad (\text{Cauchy-Schwartz}) \\
&\leq L_2 \sqrt{t} \left\| \max_{i \in [k]} X_i \right\|_{L^t}^{1/2} \left(\sqrt{\sum_{i \in [k]} \mathbb{E}[X_i]} + L_2 \left\| \sum_{i \in [k]} (X_i - \mathbb{E}[X_i]) \right\|_{L^t}^{1/2} \right).
\end{aligned}$$

Now let $C = \left\| \sum_{i \in [k]} (X_i - \mathbb{E}[X_i]) \right\|_{L^t}^{1/2}$, $B = L_2 \sqrt{\sum_{i \in [k]} \mathbb{E}[X_i]}$, and $A = \sqrt{t} \left\| \max_{i \in [k]} X_i \right\|_{L^t}^{1/2}$. then we have shown $C^2 \leq A(B + C)$. That implies C is smaller than the largest of the roots of the quadratic. Solving this quadratic inequality gives $C^2 \leq L_3(AB + A^2)$ which is the result. \square

We can now prove that SHRT and TensorSRHT has the Strong JL Moment Property.

Lemma 21. *There exists a universal constant L , such that, the following holds. Let $k \in \mathbb{Z}_{>0}$, and $(D^{(i)})_{i \in [k]} \in \prod_{i \in [k]} \mathbb{R}^{d_i \times d_i}$ be independent diagonal matrices with independent Rademacher variables. Define $d = \prod_{i \in [k]} d_i$ and $D = D_1 \times D_2 \times \dots \times D_k \in \mathbb{R}^{d \times d}$. Let $P \in \mathbb{R}^{m \times d}$ be an independent sampling matrix which samples exactly one coordinate per row, and define $M = PHD$ where H is a $d \times d$ Hadamard matrix. Let $x \in \mathbb{R}^d$ be any vector with $\|x\|_2 = 1$ and $t \geq 1$, then*

$$\left\| \frac{1}{m} \|PHDx\|_2^2 - 1 \right\|_{L^t} \leq L \left(\sqrt{\frac{tr^k}{m}} + \frac{tr^k}{m} \right),$$

where $r = \max\{t, \log m\}$.

There exists a universal constant L' , such that, setting $m = \Omega\left(\varepsilon^{-2} \log(1/\delta) (L' \log(1/\varepsilon\delta))^k\right)$, we get that $\frac{1}{\sqrt{m}} PHD$ has Strong (ε, δ) -JL Moment Property.

Note that setting $k = 1$, this matches the Fast Johnson Lindenstrauss analysis in [CNW16b].

Proof. Throughout the proof C_1, C_2 and C_3 will denote universal constants.

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

We use Lemma 20, the triangle inequality, and Cauchy-Schwartz to get that

$$\begin{aligned}
& \left\| \frac{1}{m} \sum_{i \in [m]} Z_i^2 - 1 \right\|_{L^t} \\
&= \left\| \mathbb{E} \left[\left(\frac{1}{m} \sum_{i \in [m]} Z_i^2 - 1 \right)^t \middle| D \right] \right\|_{L^t}^{1/t} \\
&\leq C_1 \left\| \frac{\sqrt{t}}{m} \mathbb{E} \left[\left(\max_{i \in [m]} Z_i^2 \right)^t \middle| D \right] \right\|_{L^t}^{1/(2t)} \sqrt{\sum_{i \in [m]} \mathbb{E}[Z_i^2 | D]} + \frac{t}{m} \mathbb{E} \left[\left(\max_{i \in [m]} Z_i^2 \right)^t \middle| D \right] \right\|_{L^t}^{1/t} \\
&\leq C_1 \frac{\sqrt{t}}{m} \left\| \mathbb{E} \left[\left(\max_{i \in [m]} Z_i^2 \right)^t \middle| D \right] \right\|_{L^t}^{1/(2t)} \left\| \sqrt{\sum_{i \in [m]} \mathbb{E}[Z_i^2 | D]} \right\|_{L^t} + C_1 \frac{t}{m} \left\| \max_{i \in [m]} Z_i^2 \right\|_{L^t} \\
&\leq C_1 \frac{\sqrt{t}}{m} \left\| \max_{i \in [m]} Z_i^2 \right\|_{L^t}^{1/2} \left\| \sum_{i \in [m]} \mathbb{E}[Z_i^2 | D] \right\|_{L^t}^{1/2} + C_1 \frac{t}{m} \left\| \max_{i \in [m]} Z_i^2 \right\|_{L^t}.
\end{aligned}$$

By orthogonality of H we have $\|HDx\|_2^2 = d\|x\|_2^2$ independent of D . Hence

$$\sum_{i \in [m]} \mathbb{E}[Z_i^2 | D] = \sum_{i \in [m]} \|x\|_2^2 = m.$$

To bound $\left\| \max_{i \in [m]} Z_i^2 \right\|_{L^t}$ we first use Lemma 19 to show

$$\left\| Z_i^2 \right\|_{L^r} = \|H_{P_i} Dx\|_{L^{2r}}^2 = \|Dx\|_{L^{2r}}^2 \leq r^k \|x\|_2^2.$$

We then bound the maximum using a sufficiently high powered sum:

$$\left\| \max_{i \in [m]} Z_i^2 \right\|_{L^t} \leq \left\| \max_{i \in [m]} Z_i^2 \right\|_{L^r} \leq \left(\sum_{i \in [m]} \left\| Z_i^2 \right\|_{L^r}^r \right)^{1/r} \leq m^{1/r} r^k \|x\|_2^2 \leq e r^k,$$

where the last inequality follows from $r \geq \log m$. This gives us that

$$\left\| \frac{1}{m} \sum_{i \in [m]} Z_i^2 - \|x\|_2^2 \right\|_{L^t} \leq C_2 \sqrt{\frac{tr^k}{m}} + C_2 \frac{tr^k}{m},$$

which finishes the first part of the proof.

We set $m = 4e^2 C_2^2 \varepsilon^{-2} \log(1/\delta) (C_3 \log(1/(\delta\varepsilon)))^k$, such that, $r \leq C_3 \log(1/(\delta\varepsilon))$. Hence $m \geq 4e^2 C_2^2 \varepsilon^{-2} \log(1/\delta) r^k$. We then get that

$$\left\| \|PHDx\|_2^2 - 1 \right\|_{L^t} \leq C_2 \sqrt{\frac{tr^k}{4e^2 C_2^2 \varepsilon^{-2} \log(1/\delta) r^k}} + C_2 \frac{tr^k}{4e^2 C_2^2 \varepsilon^{-2} \log(1/\delta) r^k} \leq \frac{\varepsilon}{e} \sqrt{\frac{t}{\log 1/\delta}},$$

for all $1 \leq t \leq \log(1/\delta)$ which finishes the proof. \square

Now we have proved that the Strong JL Moment Property is satisfied by the SRHT, the TensorSRHT as well as OSNAP transform, but we still need to prove the usefulness of the property. Our next result remedies this and show that the Strong JL Moment Property is preserved under multiplication. We will use the following decoupling lemma which first appeared in [Hit94], but the following is roughly taken from [DIPG12], which we also recommend for readers interested in more general versions.

Lemma 22 (General decoupling, [DIPG12] Theorem 7.3.1, paraphrasing). *There exists an universal constant C_0 , such that, given any two sequences $(X_i)_{i \in [n]}$ and $(Y_i)_{i \in [n]}$ of random variables, satisfying*

1. $\Pr[Y_i > t \mid (X_j)_{j \in [i-1]}] = \Pr[X_i > t \mid (X_j)_{j \in [i-1]}]$ for every $t \in \mathbb{R}$ and for every $i \in [n]$.
2. The sequence $(Y_i)_{i \in [n]}$ is conditionally independent given $(X_i)_{i \in [n]}$.
3. $\Pr[Y_i > t \mid (X_j)_{j \in [i-1]}] = \Pr[Y_i > t \mid (X_j)_{j \in [n]}]$ for every $t \in \mathbb{R}$ and for every $i \in [n]$.

Then for all $t \geq 1$,

$$\left\| \sum_{i \in [n]} X_i \right\|_{L^t} \leq C_0 \left\| \sum_{i \in [n]} Y_i \right\|_{L^t}$$

We are now ready to state and prove the main lemma of this section. This basically says that if you take k independent JL transforms, that all have the Strong $(\varepsilon/\sqrt{k}, \delta)$ -JL Moment Property, SJLMP, then the result has the (ε, δ) SJLMP. A simple union bound would give the same result where each matrix has the $(\varepsilon/k, \delta)$ SJLMP, but that would ultimately result in a higher dependency on the tensoring dimension. A simple change to the proof shows that we only need the i th JL transform to have the $(\varepsilon/\sqrt{i}, \delta)$ SJLMP, but that ultimately makes no difference for our construction.

Lemma 23. *There exists a universal constant L , such that, for any constants $\varepsilon, \delta \in [0, 1]$ and positive integer $k \in \mathbb{Z}_{>0}$. If $M^{(1)} \in \mathbb{R}^{d_2 \times d_1}, \dots, M^{(k)} \in \mathbb{R}^{d_{k+1} \times d_c}$ are independent random matrices with the Strong $(\varepsilon/(L\sqrt{k}), \delta)$ -JL Moment Property, then the matrix $M = M^{(k)} \cdot \dots \cdot M^{(1)}$ has the Strong (ε, δ) -JL Moment Property.*

Proof. Let $x \in \mathbb{R}_1^d$ be an arbitrary, fixed unit vector, and fix $1 < t \leq \log(1/\delta)$. We define $X_i = \|M^{(i)} \cdot \dots \cdot M^{(1)}x\|_2^2$ and $Y_i = X_i - X_{i-1}$ for every $i \in [k]$. By telescoping we then have that $X_i - 1 = \sum_{j \in [i]} Y_j$. We let $(T^{(i)})_{i \in [k]}$ be independent copies of $(M^{(i)})_{i \in [k]}$ and define

$$Z_i = \|T^{(i)} \cdot M^{(i-1)} \cdot \dots \cdot M^{(1)}x\|_2^2 - \|M^{(i-1)} \cdot \dots \cdot M^{(1)}x\|_2^2,$$

for every $i \in [k]$. We get the following three properties:

1. $\Pr[Z_i > t \mid (M^{(j)})_{j \in [i-1]}] = \Pr[Y_i > t \mid (M^{(j)})_{j \in [i-1]}]$ for every $t \in \mathbb{R}$ and every $i \in [k]$.
2. The sequence $(Z_i)_{i \in [k]}$ is conditionally independent given $(M^{(i)})_{i \in [k]}$.
3. $\Pr[Z_i > t \mid (M^{(j)})_{j \in [i-1]}] = \Pr[Z_i > t \mid (M^{(j)})_{j \in [k]}]$ for every $t \in \mathbb{R}$ and for every $i \in [k]$.

This means we can use Lemma 22 to get

$$\left\| \sum_{j \in [i]} Y_j \right\|_{L^t} \leq C_0 \left\| \sum_{j \in [i]} Z_j \right\|_{L^t} . \quad (17)$$

for every $i \in [k]$.

We will prove by induction on $i \in [k]$ that

$$\|X_i - 1\|_{L^t} \leq \frac{\varepsilon}{e} \sqrt{\frac{t}{\log(1/\delta)}} \leq 1 . \quad (18)$$

For $i = 1$ we use that $M^{(1)}$ has the Strong $(\varepsilon/(L\sqrt{k}), \delta)$ -JL Moment Property and get that

$$\left\| \|M^{(1)}x\|_2^2 - 1 \right\|_{L^t} \leq \frac{\varepsilon}{eL\sqrt{k}} \sqrt{\frac{t}{\log(1/\delta)}} \leq \frac{\varepsilon}{e} \sqrt{\frac{t}{\log(1/\delta)}} .$$

Now assume that (18) is true for $i - 1$. Using (17) we get that $\|X_i - 1\|_{L^t} = \left\| \sum_{j \in [i]} Y_j \right\|_{L^t} \leq C_0 \left\| \sum_{j \in [i]} Z_j \right\|_{L^t}$. By using that $(T^{(j)})_{j \in [i]}$ has the Strong $(\varepsilon/(L\sqrt{k}), \delta)$ -JL Moment Property together with Khintchine's inequality (Lemma 17), we get that

$$\begin{aligned} \left\| \sum_{j \in [i]} Z_j \right\|_{L^t} &= \left\| \mathbb{E} \left[\left(\sum_{j \in [i]} Z_j \right)^t \mid (M^{(j)})_{j \in [i]} \right] \right\|_{L^t}^{1/t} \\ &\leq C_1 \left\| \frac{\varepsilon}{eL\sqrt{k}} \sqrt{\frac{t}{\log(1/\delta)}} \sqrt{\sum_{j \in [i]} X_j^2} \right\|_{L^t} \\ &= C_1 \frac{\varepsilon}{e} \sqrt{\frac{t}{\log(1/\delta)}} \cdot \frac{1}{L\sqrt{k}} \sqrt{\left\| \sum_{j \in [i]} X_j^2 \right\|_{L^{t/2}}} \\ &\leq C_1 \frac{\varepsilon}{e} \sqrt{\frac{t}{\log(1/\delta)}} \cdot \frac{1}{L\sqrt{k}} \sqrt{\sum_{j \in [i]} \|X_j\|_{L^t}^2} , \end{aligned}$$

where the last inequality follows from the triangle inequality. Using the triangle inequality and (18) we get that

$$\|X_j\|_{L^t} \leq 1 + \|X_j - 1\|_{L^t} \leq 2 ,$$

for every $j \in [i]$. Setting $L = 2C_0C_1$ we get that

$$\left\| \sum_{j \in [i]} Y_j \right\|_{L^t} \leq \frac{\varepsilon}{e} \sqrt{\frac{t}{\log(1/\delta)}} \cdot \frac{C_0C_1}{L\sqrt{k}} \sqrt{\sum_{j \in [i]} \|X_j\|_{L^t}^2} \quad (19)$$

$$\leq \frac{\varepsilon}{e} \sqrt{\frac{t}{\log(1/\delta)}} \cdot \frac{C_0C_1}{L\sqrt{k}} \cdot 2\sqrt{i} \quad (20)$$

$$\leq \frac{\varepsilon}{e} \sqrt{\frac{t}{\log(1/\delta)}} , \quad (21)$$

which finishes the induction. Now we have that $\left\| \|Mx\|_2^2 - 1 \right\|_{L^t} \leq \frac{\varepsilon}{e} \sqrt{\frac{t}{\log(1/\delta)}}$ so we conclude that M has Strong (ε, δ) -JL Moment Property. \square

A simple corollary of this result is a sufficient condition for our recursive sketch Π^q to have the Strong JL Moment Property.

Corollary 24 (Strong JL Moment Property for Π^q). *For any integer q which is a power of two, let $\Pi^q : \mathbb{R}^{m^q} \rightarrow \mathbb{R}^m$ be defined as in Definition 11, where both of the common distributions $S_{\text{base}} : \mathbb{R}^{m^2} \rightarrow \mathbb{R}^m$ and $T_{\text{base}} : \mathbb{R}^d \rightarrow \mathbb{R}^m$, satisfy the Strong $\left(O\left(\frac{\varepsilon}{\sqrt{q}}\right), \delta\right)$ -JL Moment Property. Then it follows that Π^q satisfies the Strong (ε, δ) -JL Moment Property.*

Proof. The proof follows from using Lemma 12 and Lemma 23. □

We conclude this section by proving Theorem 2.

Theorem 2. *For every positive integers n, p, d , every $\varepsilon, s_\lambda > 0$, there exists a distribution on linear sketches $\Pi^p \in \mathbb{R}^{m \times d^p}$ such that: (1) If $m = \tilde{\Omega}(ps_\lambda^2 \varepsilon^{-2})$, then Π^p is an $(\varepsilon, 1/\text{poly}(n), s_\lambda, d^p, n)$ -oblivious subspace embedding (Definition 2). (2) If $m = \tilde{\Omega}(p\varepsilon^{-2})$, then Π^p has the $(\varepsilon, 1/\text{poly}(n))$ -approximate matrix product property (Definition 3).*

Moreover, in the setting of (1), for any $X \in \mathbb{R}^{d \times n}$, if $A \in \mathbb{R}^{d^p \times n}$ is the matrix whose columns are obtained by a p -fold self-tensoring of each column of X , then the matrix $\Pi^p A$ can be computed using Algorithm 1 in time $\tilde{O}(pnm + p^{3/2}s_\lambda \varepsilon^{-1} \text{nnz}(X))$.

Proof. Let $\delta = \frac{1}{\text{poly}(n)}$ denote the failure probability. Define $q = \lceil \log_2(p) \rceil$ and let $\Pi^p \in \mathbb{R}^{m \times d^p}$ and $\Pi^q \in \mathbb{R}^{m \times d^q}$ be the sketches defined in Definition 11, where $S_{\text{base}} \in \mathbb{R}^{m \times m^2}$ is a TensorSRHT sketch and $T_{\text{base}} \in \mathbb{R}^{m \times d}$ is an OSNAP sketch with sparsity parameter s , which will be set later.

Oblivious Subspace Embedding Let $m = \Theta\left(\frac{ps_\lambda^2 \log(1/(\varepsilon\delta))^3}{\varepsilon^2}\right)$ and $s = \Theta\left(\frac{\sqrt{p}s_\lambda \log(1/\delta)}{\varepsilon}\right)$ be integers, then Lemma 21 and Lemma 18 implies that S_{base} and T_{base} has the Strong $\left(O\left(\frac{\varepsilon}{\sqrt{q}s_\lambda}\right), \delta\right)$ -JL Moment Property, thus using Corollary 24 we conclude that Π^q has the Strong $\left(\frac{\varepsilon}{s_\lambda}, \delta\right)$ -JL Moment Property and in particular it has the $\left(\frac{\varepsilon}{s_\lambda}, \delta, \log(1/\delta)\right)$ -JL Moment Property. By Lemma 11 we then get that Π^q is an $(\varepsilon, \delta, s_\lambda, d^q, n)$ -Oblivious Subspace Embedding, and by Lemma 10 we get that Π^p is an $(\varepsilon, \delta, s_\lambda, d^p, n)$ -Oblivious Subspace Embedding.

Approximate Matrix Multiplication Let $m = \Theta\left(\frac{p \log(1/(\varepsilon\delta))^3}{\varepsilon^2}\right)$ and $s = \Theta\left(\frac{\sqrt{p} \log(1/\delta)}{\varepsilon}\right)$ be integers. Then Lemma 21 and Lemma 18 implies that S_{base} and T_{base} has the Strong $\left(O\left(\frac{\varepsilon}{\sqrt{q}s_\lambda}\right), \delta\right)$ -JL Moment Property. Thus, using Corollary 24 we conclude that Π^q has the Strong (ε, δ) -JL Moment Property and in particular it has the $(\varepsilon, \delta, \log(1/\delta))$ -JL Moment Property. By Lemma 11 we then get that Π^q has the (ε, δ) -Approximate Matrix Multiplication Property, and by Lemma 10 we get that Π^p has the (ε, δ) -Approximate Matrix Multiplication Property.

Runtime of Algorithm 1 when the base sketch S_{base} is a TensorSRHT sketch and T_{base} is an OSNAP sketch with sparsity parameter s : We compute the time of running Algorithm 1 on a vector x . Computing Y_j^0 for each j in lines 3 and 4 of algorithm requires applying an OSNAP sketch on either x or e_1 which takes time $O(s \cdot \text{nnz}(x))$. Therefore computing all Y_j^0 's takes time $O(qs \cdot \text{nnz}(x))$.

Computing each of Y_j^l 's for $l \geq 1$ in line 7 of Algorithm 1 amounts to applying a TensorSRHT sketch of input dimension m^2 and target dimension of m on $Y_{2j-1}^{l-1} \otimes Y_{2j}^{l-1}$. This takes time

$O(m \log m)$. Therefore computing Y_j^l for all $l, j \geq 1$ takes time $O(q \cdot m \log m)$. Note that $q \leq 2p$ hence the total running time of Algorithm 1 on one vector x is $O(pm \log_2 m + ps \cdot \text{nnz}(w))$. Sketching n columns of a matrix $X \in \mathbb{R}^{d \times n}$ takes time $O(p(nm \log_2 m + s \cdot \text{nnz}(X)))$.

In the setting of (1) we have that $s = O\left(\frac{\sqrt{ps_\lambda \log(1/\delta)}}{\epsilon}\right)$, hence we get a runtime of $O\left(pnm \log_2 m + \frac{p^{3/2}s_\lambda \log(1/\delta)}{\epsilon} \text{nnz}(X)\right) = \tilde{O}\left(pnm + \frac{p^{3/2}s_\lambda}{\epsilon} \text{nnz}(X)\right)$. \square

5 Linear Dependence on the Statistical Dimension s_λ

In this section, we show that if one chooses the internal nodes and the leaves of our recursive construction from Section 3 to be TensorSRHT and OSNAP transform respectively, then the recursive construction Π^q as in Definition 11 yields a high probability OSE with target dimension $\tilde{O}(p^4 s_\lambda)$. Thus, we prove Theorem 3. This sketch is very efficiently computable for high degree tensor products because the OSNAP transform is computable in input sparsity time and the TensorSRHT supports fast matrix vector multiplication for tensor inputs.

We start by defining the *Spectral Property* for a sketch. We use the notation $\|\cdot\|_{op}$ to denote the operator norm of matrices.

Definition 20 (Spectral Property). For any positive integers m, n, d and any $\epsilon, \delta, \mu_F, \mu_2 \geq 0$ we say that a random matrix $S \in \mathbb{R}^{m \times d}$ satisfies the $(\mu_F, \mu_2, \epsilon, \delta, n)$ -spectral property if, for every fixed matrix $U \in \mathbb{R}^{d \times n}$ with $\|U\|_F^2 \leq \mu_F$ and $\|U\|_{op}^2 \leq \mu_2$,

$$\Pr_S \left[\left\| U^\top S^\top S U - U^\top U \right\|_{op} \leq \epsilon \right] \geq 1 - \delta.$$

The *spectral property* is a central property of our sketch construction from Section 3 when leaves are OSNAP and internal nodes are TensorSRHT. This is a powerful property which implies that any sketch which satisfies the *spectral property*, is an *Oblivious Subspace Embedding*. The SRHT, TensorSRHT, as well as OSNAP sketches (Definitions 14, 15, 16 respectively) with target dimension $m = \Omega\left(\left(\frac{\mu_F \mu_2}{\epsilon^2}\right) \cdot \text{poly}(\log(nd/\delta))\right)$ and sparsity parameter $s = \Omega(\text{poly}(\log(nd/\delta)))$, all satisfy the above-mentioned spectral property [Sar06, Tro11, NN13].

In section 5.1 we recall the tools from the literature which we use to prove the spectral property for our construction Π^q . Then in section 5.2 we show that our recursive construction in Section 3 satisfies the Spectral Property of Definition 20 as long as $I_{d^q} \times T_{\text{base}}$ and $I_{m^q} \times S_{\text{base}}$ satisfy the Spectral Property. Therefore, we analyze the Spectral Property of $I_{d^q} \times \text{OSNAP}$ and $I_{m^q} \times \text{TensorSRHT}$ in section 5.3 and section 5.4 respectively. Finally we put everything together in section 5.5 and prove that when the leaves are OSNAP and the internal nodes are TensorSRHT in our recursive construction of Section 3, the resulting sketch Π^q satisfies the Spectral Property thereby proving Theorem 3.

5.1 Matrix Concentration Tools

In this section we present the definitions and tools which we use for proving concentration properties of random matrices.

Claim 25. For every $\epsilon, \delta > 0$ and any sketch $S \in \mathbb{R}^{m \times d}$ such that $I_k \times S$ satisfies $(\mu_F, \mu_2, \epsilon, \delta, n)$ -spectral property, the sketch $S \times I_k$ also satisfies the $(\mu_F, \mu_2, \epsilon, \delta, n)$ -spectral property.

Proof. Suppose $U \in \mathbb{R}^{dk \times n}$. Then, note that there exists $U' \in \mathbb{R}^{dk \times n}$ formed by permuting the rows of U such that $(S \times I_k)U$ and $(I_k \times S)U'$ are identical up to a permutation of the rows. (In

particular, U' is the matrix such that the (d, k) -reshaping of any column U^j of U' is the transpose of the (k, d) -reshaping of the corresponding column U'^j of U' .) Then, observe that

$$U^\top U = U'^\top U'.$$

and

$$U^\top (S \times I_k)^\top (S \times I_k) U = U'^\top (I_k \times S)^\top (I_k \times S) U'.$$

Therefore,

$$\|U^\top (S \times I_k)^\top (S \times I_k) U - U^\top U\|_{op} = \|U'^\top (S \times I_k)^\top (S \times I_k) U' - U'^\top U'\|_{op}.$$

Moreover, since U and U' are identical up to a permutation of the rows, we have $\|U\|_{op} = \|U'\|_{op}$ and $\|U\|_F = \|U'\|_F$. The desired claim now follows easily. \square

We will use matrix Bernstein inequalities to show spectral guarantees for sketches,

Lemma 26 (Matrix Bernstein Inequality (Theorem 6.1.1 in [Tro15])). *Consider a finite sequence Z_i of independent, random matrices with dimensions $d_1 \times d_2$. Assume that each random matrix satisfies $\mathbb{E}[Z_i] = 0$ and $\|Z_i\|_{op} \leq B$ almost surely. Define $\sigma^2 = \max\{\|\sum_i \mathbb{E}[Z_i Z_i^*]\|_{op}, \|\sum_i \mathbb{E}[Z_i^* Z_i]\|_{op}\}$. Then for all $t > -0$,*

$$\mathbb{P}\left[\left\|\sum_i Z_i\right\|_{op} \geq t\right] \leq (d_1 + d_2) \cdot \exp\left(\frac{-t^2/2}{\sigma^2 + Bt/3}\right).$$

Lemma 27 (Restatement of Corollary 6.2.1 of [Tro15]). *Let B be a fixed $n \times n$ matrix. Construct an $n \times n$ matrix R that satisfies,*

$$\mathbb{E}[R] = B \quad \text{and} \quad \|R\|_{op} \leq L,$$

almost surely. Define $M = \max\{\|\mathbb{E}[RR^]\|_{op}, \|\mathbb{E}[R^*R]\|_{op}\}$. Form the matrix sampling estimator,*

$$\bar{R} = \frac{1}{m} \sum_{k=1}^m R_k,$$

where each R_k is an independent copy of R . Then,

$$\Pr\left[\|\bar{R} - B\|_{op} \geq \epsilon\right] \leq 8n \cdot \exp\left(\frac{-m\epsilon^2/2}{M + 2L\epsilon/3}\right).$$

To analyze the performance of SRHT we need the following claim which shows that with high probability individual entries of the Hadamard transform of a vector with random signs on its entries do not “overshoot the mean energy” by much.

Claim 28. *Let D be a $d \times d$ diagonal matrix with independent Rademacher random variables along the diagonal. Also, let H be a $d \times d$ Hadamard matrix. Then, for every $x \in \mathbb{R}^d$,*

$$\Pr\left[\|HD \cdot x\|_\infty \leq 2\sqrt{\log_2(d/\delta)} \cdot \|x\|_2\right] \geq 1 - \delta.$$

Proof. By Khintchine's inequality, Lemma 17 we have that for every $t \geq 1$ and every $j \in [d]$ the j^{th} element of HDx has a bounded t^{th} moment as follows,

$$\|(HDx)_j\|_{L^t} \leq \sqrt{t} \cdot \|x\|_2.$$

Hence by applying Markov's inequality to the t^{th} moment of $|(HDx)_j|$ for $t = \log_2(d/\delta)$ we get that,

$$\Pr \left[|(HDx)_j| \geq 2\sqrt{\log_2(d/\delta)} \cdot \|x\|_2 \right] \leq \delta/d.$$

The claim follows by a union bound over all entries $j \in [d]$. \square

Claim 29. *Let D_1, D_2 be two independent $d \times d$ diagonal matrices, each with diagonal entries given by independent Rademacher random variables. Also, let H be a $d \times d$ Hadamard matrix. Then, for every $x \in \mathbb{R}^{d^2}$,*

$$\Pr_{D_1, D_2} \left[\|(HD_1) \times (HD_2) \cdot x\|_\infty \leq 4\log_2(d/\delta) \cdot \|x\|_2 \right] \geq 1 - \delta.$$

Proof. By Claim 6 we can write that,

$$(HD_1) \times (HD_2) = (H \times H)(D_1 \times D_2),$$

where $H \times H$ is indeed a Hadamard matrix of size $d^2 \times d^2$ which we denote by H' . The goal is to prove

$$\Pr_{D_1, D_2} \left[\|H'(D_1 \times D_2) \cdot x\|_\infty \leq 4\log_2(d/\delta) \cdot \|x\|_2 \right] \geq 1 - \delta.$$

By Lemma 19 we have that for every $t \geq 1$ and every $j \in [d^2]$ the j^{th} element of $H'(D_1 \times D_2)x$ has a bounded t^{th} moment as follows,

$$\|(H'(D_1 \times D_2)x)_j\|_{L^t} \leq t \cdot \|x\|_2.$$

Hence by applying Markov's inequality to the t^{th} moment of $|(H'(D_1 \times D_2)x)_j|$ for $t = \log_2(d/\delta)$ we get that,

$$\Pr \left[|(H'(D_1 \times D_2)x)_j| \geq 4\log_2(d/\delta) \cdot \|x\|_2 \right] \leq \delta/d^2.$$

The claim follows by a union bound over all entries $j \in [d^2]$. \square

5.2 Spectral Property of the sketch Π^q

In this section we show that the sketch Π^q presented in Definition 11 inherits the spectral property (see Definition 20) from the base sketches S_{base} and T_{base} . We start by the following claim which proves that composing two random matrices with spectral property results in a matrix with spectral property.

Claim 30. *For every $\epsilon, \epsilon', \delta, \delta' > 0$, suppose that $S \in \mathbb{R}^{m \times t}$ is a sketch which satisfies the $((\mu_F + 1)(1 + \epsilon'), \mu_2 + 1 + \epsilon', \epsilon, \delta, n)$ -spectral property and also suppose that the sketch $T \in \mathbb{R}^{t \times d}$ satisfies the $(\mu_F + 1, \mu_2 + 1, \epsilon', \delta'/n, n)$ -spectral property. Then $S \cdot T$ satisfies the $(\mu_F + 1, \mu_2 + 1, \epsilon + \epsilon', \delta + \delta'(1 + 1/n), n)$ -spectral property.*

Proof. Suppose S and T are matrices satisfying the hypothesis of the claim. Consider an arbitrary matrix $U \in \mathbb{R}^{d \times n}$ which satisfies $\|U\|_F^2 \leq \mu_F + 1$ and $\|U\|_{op}^2 \leq \mu_2 + 1$. We want to prove that for every such U ,

$$\Pr \left[\|U^\top (S \cdot T)^\top (S \cdot T) U - U^\top U\|_{op} \leq \epsilon + \epsilon' \right] \geq 1 - \delta - \delta'(1 + 1/n).$$

Let us define the event \mathcal{E} as follows,

$$\mathcal{E} := \left\{ \|T \cdot U\|_F^2 \leq (1 + \epsilon') \|U\|_F^2 \text{ and } \|U^\top T^\top T U - U^\top U\|_{op} \leq \epsilon' \right\}.$$

We show that this event holds with probability $1 - \delta'(1 + 1/n)$ over the random choice of sketch T . The spectral property of T implies that for every column U^j of matrix U ,

$$\|TU^j\|_2^2 = (1 \pm \epsilon') \|U^j\|_2^2,$$

with probability $1 - \frac{\delta'}{n}$. By a union bound over all $j \in [n]$, we have the following,

$$\Pr_T \left[\|T \cdot U\|_F^2 \leq (1 + \epsilon') \|U\|_F^2 \right] \geq 1 - \delta'.$$

Also,

$$\Pr_T \left[\|U^\top T^\top T U - U^\top U\|_{op} \leq \epsilon' \right] \geq 1 - \delta'/n.$$

Therefore by union bound,

$$\Pr_T[\mathcal{E}] \geq 1 - \delta'(1 + 1/n).$$

We condition on $T \in \mathcal{E}$ in the rest of the proof. Since S satisfies the $((\mu_F + 1)(1 + \epsilon'), \mu_2 + 1 + \epsilon', \epsilon, \delta, n)$ -spectral property,

$$\Pr_S \left[\|(TU)^\top S^\top S(TU) - (TU)^\top (TU)\|_{op} \leq \epsilon \right] \geq 1 - \delta.$$

Therefore,

$$\begin{aligned} & \Pr_{T,S} \left[\|U^\top (S \cdot T)^\top (S \cdot T) U - U^\top U\|_{op} \leq \epsilon + \epsilon' \right] \\ & \geq \Pr_S \left[\|U^\top (S \cdot T)^\top (S \cdot T) U - U^\top U\|_{op} \leq \epsilon + \epsilon' \mid T \in \mathcal{E} \right] - \Pr_T[\bar{\mathcal{E}}] \\ & \geq \Pr_S \left[\|(TU)^\top S^\top S(TU) - U^\top U\|_{op} \leq \epsilon + \epsilon' \mid T \in \mathcal{E} \right] - \delta'(1 + 1/n) \\ & \geq \Pr_S \left[\|(TU)^\top S^\top S(TU) - (TU)^\top (TU)\|_{op} + \|(TU)^\top (TU) - U^\top U\|_{op} \leq \epsilon + \epsilon' \mid T \in \mathcal{E} \right] - \delta'(1 + \frac{1}{n}) \\ & \geq \Pr_S \left[\|(TU)^\top S^\top S(TU) - (TU)^\top (TU)\|_{op} \leq \epsilon \mid T \in \mathcal{E} \right] - \delta'(1 + 1/n) \\ & \geq 1 - \delta - \delta'(1 + 1/n). \end{aligned}$$

This completes the proof. \square

In the following lemma we show that composing independent random matrices with spectral property preserves the spectral property.

Lemma 31. For any $\varepsilon, \delta, \mu_F, \mu_2 > 0$ and every positive integers k, n , if $M^{(1)} \in \mathbb{R}^{d_2 \times d_1}, \dots, M^{(k)} \in \mathbb{R}^{d_{k+1} \times d_k}$ are independent random matrices with the $(2\mu_F + 2, 2\mu_2 + 2, O(\varepsilon/k), O(\delta/nk), n)$ -spectral property then the product matrix $M = M^{(k)} \dots M^{(1)}$ satisfies the $(\mu_F + 1, \mu_2 + 1, \varepsilon, \delta, n)$ -spectral property.

Proof. Consider a matrix $U \in \mathbb{R}^{d_1 \times n}$ which satisfies $\|U\|_F^2 \leq \mu_F + 1$ and $\|U\|_{op}^2 \leq \mu_2 + 1$. We want to prove that for every such U ,

$$\Pr \left[\|U^\top M^\top M U - U^\top U\|_{op} \leq \varepsilon \right] \geq 1 - \delta,$$

where $M = M^{(k)} \dots M^{(1)}$.

By the assumption of the lemma the matrices $M^{(1)}, \dots, M^{(k)}$ satisfy the $(2\mu_F + 2, 2\mu_2 + 2, O(\varepsilon/k), O(\delta/nk), n)$ -spectral property. For every $j \in [k]$, let us define the set \mathcal{E}_j as follows,

$$\mathcal{E}_j := \left\{ \left(M^{(1)}, \dots, M^{(j)} \right) : \left\{ \begin{array}{l} 1. \left\| \left(M^{(j)} \dots M^{(1)} \right) U \right\|_F^2 \leq \left(1 + \frac{\varepsilon}{10k} \right)^j \|U\|_F^2 \\ 2. \left\| U^\top \left(M^{(j)} \dots M^{(1)} \right)^\top \left(M^{(j)} \dots M^{(1)} \right) U - U^\top U \right\|_{op} \leq \frac{\varepsilon_j}{3k} \end{array} \right\}.$$

First we prove that for every $j \in \{1, \dots, k-1\}$,

$$\Pr_{M^{(j+1)}} \left[\left(M^{(1)}, \dots, M^{(j+1)} \right) \in \mathcal{E}_{j+1} \mid \left(M^{(1)}, \dots, M^{(j)} \right) \in \mathcal{E}_j \right] \geq 1 - \frac{\delta}{2k}.$$

Let us denote $\left(M^{(j)} \dots M^{(1)} \right) \cdot U$ by U' . The condition $\left(M^{(1)}, \dots, M^{(j)} \right) \in \mathcal{E}_j$ implies that, $\|U'\|_F^2 \leq \left(1 + \varepsilon/(10k) \right)^j \|U\|_F^2$ and $\|U'^\top U' - U^\top U\|_{op} \leq \frac{\varepsilon_j}{3k}$ and therefore by triangle inequality we have $\|U'\|_{op}^2 \leq \left(\|U\|_{op} + \frac{\varepsilon_j}{3k} \right)^2$. The assumptions $\|U\|_F^2 \leq \mu_F + 1$ and $\|U\|_{op}^2 \leq \mu_2 + 1$ imply that $\|U'\|_F^2 \leq 2\mu_F + 2$ and $\|U'\|_{op}^2 \leq 2\mu_2 + 2$. Now note that by the assumption of the lemma, $M^{(j+1)}$ satisfies the $(2\mu_F + 2, 2\mu_2 + 2, O(\varepsilon/k), O(\delta/nk), n)$ -spectral property. Therefore,

$$\Pr_{M^{(j+1)}} \left[\left\| \left(M^{(j+1)} U' \right)^\top M^{(j+1)} U' - U'^\top U' \right\|_{op} \leq \frac{\varepsilon}{3k} \mid \left(M^{(1)}, \dots, M^{(j)} \right) \in \mathcal{E}_j \right] \geq 1 - \delta/(4nk).$$

Combining the above with $\|U'^\top U' - U^\top U\|_2 \leq \frac{\varepsilon_j}{3k}$ gives,

$$\Pr_{M^{(j+1)}} \left[\left\| \left(M^{(j+1)} U' \right)^\top M^{(j+1)} U' - U^\top U \right\|_{op} \leq \varepsilon \frac{j+1}{3k} \mid \left(M^{(1)}, \dots, M^{(j)} \right) \in \mathcal{E}_j \right] \geq 1 - \delta/(4nk). \quad (22)$$

Also from the spectral property of $M^{(j+1)}$ it follows that for every column U'^i of matrix U' ,

$$\|M^{(j+1)} U'^i\|_2^2 = (1 \pm \varepsilon/(10k)) \|U'^i\|_2^2,$$

with probability $1 - \frac{\delta}{4nk}$. By a union bound over all $i \in [n]$, we have the following,

$$\Pr_{M^{(j+1)}} \left[\|M^{(j+1)} \cdot U'\|_F^2 \leq \left(1 + \varepsilon/(10k) \right) \|U'\|_F^2 \mid \left(M^{(1)}, \dots, M^{(j)} \right) \in \mathcal{E}_j \right] \geq 1 - \frac{\delta}{4k}.$$

Combining the above with $\|U'\|_F^2 \leq \left(1 + \varepsilon/(10k) \right)^j \|U\|_F^2$ gives,

$$\Pr_{M^{(j+1)}} \left[\|M^{(j+1)} \cdot U'\|_F^2 \leq \left(1 + \frac{\varepsilon}{10k} \right)^{j+1} \|U\|_F^2 \mid \left(M^{(1)}, \dots, M^{(j)} \right) \in \mathcal{E}_j \right] \geq 1 - \frac{\delta}{4k}. \quad (23)$$

A union bound on (22) and (23) gives,

$$\Pr_{M^{(j+1)}} \left[\left(M^{(1)}, \dots, M^{(j+1)} \right) \in \mathcal{E}_{j+1} \mid \left(M^{(1)}, \dots, M^{(j)} \right) \in \mathcal{E}_j \right] \geq 1 - \frac{\delta}{4nk} - \frac{\delta}{4k} \geq 1 - \frac{\delta}{2k}.$$

We also show that,

$$\Pr_{M^{(1)}} [M^{(1)} \in \mathcal{E}_1] \geq 1 - \delta/2k.$$

By the assumption of lemma we know that $M^{(1)}$ satisfies the $\left(2\mu_F + 2, 2\mu_2 + 2, \frac{\epsilon}{10k}, \frac{\delta}{4nk}, n\right)$ -spectral property. Therefore,

$$\Pr_{M^{(1)}} \left[\|(M^{(1)}U)^\top M^{(1)}U - U^\top U\|_{op} \leq \frac{\epsilon}{10k} \right] \geq 1 - \frac{\delta}{4nk}. \quad (24)$$

Also for every column U^i of matrix U ,

$$\|M^{(1)}U^i\|_2^2 = (1 \pm \epsilon/(10k)) \|U^i\|_2^2,$$

with probability $1 - \frac{\delta}{4nk}$. By a union bound over all $i \in [n]$, we have the following,

$$\Pr_{M^{(1)}} \left[\|M^{(1)} \cdot U\|_F^2 \leq (1 + \epsilon/(10k)) \|U\|_F^2 \right] \geq 1 - \frac{\delta}{4k}. \quad (25)$$

A union bound on (24) and (25) gives,

$$\Pr_{T_1} [T_1 \in \mathcal{E}_1] \geq 1 - \frac{\delta}{4nk} - \frac{\delta}{4k} \geq 1 - \frac{\delta}{2k}.$$

By the chain rule for events we have,

$$\begin{aligned} & \Pr_{M^{(1)}, \dots, M^{(k)}} \left[\left(M^{(1)}, \dots, M^{(k)} \right) \in \mathcal{E}_k \right] \\ & \geq \prod_{j=2}^k \Pr_{M^{(j)}} \left[\left(M^{(1)}, \dots, M^{(j)} \right) \in \mathcal{E}_j \mid \left(M^{(1)}, \dots, M^{(j-1)} \right) \in \mathcal{E}_{j-1} \right] \cdot \Pr_{M^{(1)}} [M^{(1)} \in \mathcal{E}_1] \\ & \geq \left(1 - \frac{\delta}{2k}\right)^k \geq 1 - \delta, \end{aligned}$$

which completes the proof of the lemma. \square

The following lemma shows that our sketch construction Π^q presented in definition 11 inherits the spectral property of Definition 20 from the base sketches, that is, if S_{base} and T_{base} are such that $I_{m^{q-2}} \times S_{\text{base}}$ and $I_{d^{q-1}} \times T_{\text{base}}$ satisfy the spectral property, then the sketch Π^q satisfies the spectral property.

Lemma 32. *For every positive integers n, d, m , any power of two integer q , any base sketch $T_{\text{base}} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ such that $I_{d^{q-1}} \times T_{\text{base}}$ satisfies the $(2\mu_F + 2, 2\mu_2 + 2, O(\epsilon/q), O(\delta/nq), n)$ -spectral property, any $S_{\text{base}} : \mathbb{R}^{m^2} \rightarrow \mathbb{R}^m$ such that $I_{m^{q-2}} \times S_{\text{base}}$ satisfies the $(2\mu_F + 2, 2\mu_2 + 2, O(\epsilon/q), O(\delta/nq), n)$ -spectral property, the sketch Π^q defined as in Definition 11 satisfies the $(\mu_F + 1, \mu_2 + 1, \epsilon, \delta, n)$ -spectral property.*

Proof. We wish to show that $\Pi^q = Q^q T^q$ as per Definition 11, satisfies the $(\mu_F + 1, \mu_2 + 1, \epsilon, \delta, n)$ -spectral property. By Definition 9 $Q^q = S^2 S^4 \dots S^q$. Claim 7 shows that for every $l \in \{2, 4, \dots, q\}$ we can write,

$$S^l = M_{l/2}^l M_{l/2-1}^l \dots M_1^l, \quad (26)$$

where $M_j = I_{m^{q-2j}} \times S_{q/2-j+1}^q \times I_{m^{j-1}}$ for every $j \in [q/2]$. From the discussion in Definition 10 it follows that,

$$T^q = M'_q \dots M'_1, \quad (27)$$

where $M'_j = I_{d^{q-j}} \times T_{q-j+1} \times I_{m^{j-1}}$ for every $j \in [q]$. Therefore by combining (26) and (27) we get that,

$$\Pi^q = M^{(2q+1)} M^{(2q)} \dots M^{(1)},$$

where $M^{(i)}$ matrices are independent and by the assumption of the lemma about the spectral property of $I_{m^{q-2}} \times S_{\text{base}}$ and $I_{d^{q-1}} \times T_{\text{base}}$ together with Claim 25 it follows that $M^{(i)}$ matrices satisfy the $(2\mu_F + 2, 2\mu_2 + 2, O(\epsilon/q), O(\delta/nq), n)$ -spectral property. Therefore, the Lemma readily follows by invoking Lemma 31 with $k = 2q + 1$. \square

5.3 Spectral Property of Identity \times TensorSRHT

In this section, we show that tensoring an identity operator with a TensorSRHT sketch results in a transform that satisfies the spectral property defined in Definition 20 with nearly optimal target dimension.

Lemma 33. *Suppose $\epsilon, \delta, \mu_2, \mu_F > 0$ and n is a positive integer. If $m = \Omega\left(\log\left(\frac{n}{\delta}\right) \log^2\left(\frac{ndk}{\epsilon\delta}\right) \cdot \frac{\mu_F \mu_2}{\epsilon^2}\right)$ and $S \in \mathbb{R}^{m \times d}$ is a TensorSRHT, then the sketch $I_k \times S$ satisfies $(\mu_F, \mu_2, \epsilon, \delta, n)$ -spectral property.*

Proof. Fix a matrix $U \in \mathbb{R}^{kd \times n}$ with $\|U\|_F^2 \leq \mu_F$ and $\|U\|_{op}^2 \leq \mu_2$. Partition U by rows into $d \times n$ submatrices U_1, U_2, \dots, U_k such that $U^\top = \begin{bmatrix} U_1^\top & U_2^\top & \dots & U_k^\top \end{bmatrix}$. Note that

$$U^\top (I_k \times S)^\top (I_k \times S) U = (U_1)^\top S^\top S U_1 + \dots + (U_k)^\top S^\top S U_k.$$

The proof first considers the simpler case of a TensorSRHT sketch of rank 1 and then applies the matrix Bernstein inequality from Lemma 27. Let R denote a rank one TensorSRHT sketch. R is a $1 \times d$ matrix defined in Definition 15 by setting $m = 1$ as follows,

$$R = P \cdot (HD_1 \times HD_2),$$

where $P \in \{0, 1\}^{1 \times d}$ has one non-zero element whose position is uniformly distributed over $[d]$. Note that $S^\top S \in \mathbb{R}^{d \times d}$, is the average of m independent samples from $R^\top R$, i.e., $S^\top S = \frac{1}{m} \sum_{i \in [m]} R_i^\top R_i$, for i.i.d. $R_1, R_2, \dots, R_m \sim R$, and therefore,

$$U^\top (I_k \times S)^\top (I_k \times S) U = \frac{1}{m} \sum_{i \in [m]} U^\top (I_k \times R_i)^\top (I_k \times R_i) U.$$

Therefore in order to use matrix Bernstein, Lemma 27, we need to bound the maximum operator norm of $U^\top (I_k \times R)^\top (I_k \times R) U$ as well as the operator norm of its second moment.

We proceed to upper bound the operator norm of $U^\top(I_k \times R)^\top(I_k \times R)U$. First, define the set $\mathcal{E} := \left\{ (D_1, D_2) : \left\| (HD_1 \times HD_2)U_j^i \right\|_\infty^2 \leq 16\log^2\left(\frac{nd\mu_F k}{\epsilon\delta}\right) \cdot \|U_j^i\|_2^2 \text{ for all } j \in [k] \text{ and all } i \in [n] \right\}$,

where U_j^i is the i th column of U^j . By Claim 29, for every $i \in [n]$ and $j \in [k]$,

$$\Pr_{D_1, D_2} \left[\left\| (HD_1 \times HD_2)U_j^i \right\|_\infty^2 \leq 16\log^2(ndk/\delta)\|U_j^i\|_2^2 \right] \geq 1 - \epsilon\delta/(nk\mu_F d).$$

Thus, by a union bound over all $i \in [n]$ and $j \in [k]$, it follows that \mathcal{E} occurs with probability at least $1 - \epsilon\delta/(d\mu_F)$,

$$\Pr_{D_1, D_2} [(D_1, D_2) \in \mathcal{E}] \geq 1 - \epsilon\delta/(d\mu_F),$$

where the probability is over the random choice of D_1, D_2 .

From now on, we fix $(D_1, D_2) \in \mathcal{E}$ and proceed having conditioned on this event.

Upper bounding $\left\| U^\top(I_k \times R)^\top(I_k \times R)U \right\|_{op}$. From the fact that we have conditioned on $(D_1, D_2) \in \mathcal{E}$, note that

$$\begin{aligned} L \equiv \left\| U^\top(I_k \times R)^\top(I_k \times R)U \right\|_{op} &= \|(U^1)^\top R^\top R U_1 + \cdots + (U_k)^\top R^\top R U_k\|_{op} \\ &\leq \left\| (U_1)^\top R^\top R U_1 \right\|_{op} + \cdots + \left\| (U_k)^\top R^\top R U_k \right\|_{op} \\ &= \|R U_1\|_2^2 + \cdots + \|R U_k\|_2^2 \\ &\leq 16\log^2(nd\mu_F k/\epsilon\delta) \cdot (\|U_1\|_F^2 + \cdots + \|U_k\|_F^2) \\ &\leq 16\log^2(nd\mu_F k/\epsilon\delta) \cdot \|U\|_F^2 \\ &= 16\mu_F \cdot \log^2(nd\mu_F k/\epsilon\delta), \end{aligned}$$

where the equality on the third line above holds because the matrices $(U^i)^\top R^\top R U^i$ are rank one.

Upper bounding $\left\| \mathbb{E}_P \left[\left(U^\top(I_k \times R)^\top(I_k \times R)U \right)^2 \right] \right\|_{op}$. For every $x \in \mathbb{R}^d$ with $\|x\|_2 = 1$, we have

$$\begin{aligned} x^T \mathbb{E}_P \left[\left(U^\top(I_k \times R)^\top(I_k \times R)U \right)^2 \right] x &= \mathbb{E}_P \left[\sum_{j, j' \in [k]} x^T (U_j)^\top R^\top R U_j \cdot (U_{j'})^\top R^\top R U_{j'} x \right] \\ &\leq \mathbb{E}_P \left[\sum_{j, j' \in [k]} |R U_j x| \|R U_j\|_2 |R U_{j'} x| \|R U_{j'}\|_2 \right] \\ &= \mathbb{E}_P \left[\left(\sum_{j \in [k]} |R U_j x| \|R U_j\|_2 \right)^2 \right] \\ &\leq \mathbb{E}_P \left[\left(\sum_{j \in [k]} (R U_j x)^2 \right) \left(\sum_{j \in [k]} \|R U_j\|_2^2 \right) \right], \end{aligned}$$

where the second and fourth lines follow from the Cauchy-Schwarz inequality. Using the fact that we conditioned on $(D_1, D_2) \in \mathcal{E}$, we get

$$\begin{aligned}
x^T \mathbb{E}_P \left[\left(U^\top (I_k \times R)^\top (I_k \times R) U \right)^2 \right] x &\leq 16 \log^2(nd\mu_F k/\epsilon\delta) \left(\sum_{j \in [k]} \|U_j\|_F^2 \right) \mathbb{E}_P \left[\sum_{j \in [k]} (RU_j x)^2 \right] \\
&= 16 \log^2(nd\mu_F k/\epsilon\delta) \left(\sum_{j \in [k]} \|U_j\|_F^2 \right) \sum_{j \in [k]} \mathbb{E}_P \left[(P(HD_1 \times HD_2)U_j x)^2 \right] \\
&= 16 \log^2(nd\mu_F k/\epsilon\delta) \cdot \|U\|_F^2 \sum_{j \in [k]} \|U_j x\|_2^2 \\
&= 16 \log^2(nd\mu_F k/\epsilon\delta) \cdot \|U\|_F^2 \|Ux\|_2^2 \\
&\leq 16 \log^2(nd\mu_F k/\epsilon\delta) \cdot \mu_F \mu_2,
\end{aligned}$$

since $\mathbb{E}_P [(P(HD_1 \times HD_2)U_j x)^2] = \frac{1}{d} \|(HD_1 \times HD_2)U_j x\|^2 = \|U_j x\|_2^2$ for all x .

Since the matrix $\mathbb{E}_P \left[\left(U^\top (I_k \times R)^\top (I_k \times R) U \right)^2 \right]$ is positive semi-definite for any fixed D_1 and D_2 , it follows that

$$M \equiv \left\| \mathbb{E}_P \left[\left(U^\top (I_k \times R)^\top (I_k \times R) U \right)^2 \right] \right\|_{op} \leq 16 \log^2(nd\mu_F k/\epsilon\delta) \cdot \mu_F \mu_2.$$

Combining one-dimensional TensorSRHT sketches. To conclude, we note that the Gram matrix of a TensorSRHT, $S^\top S \in \mathbb{R}^{d \times d}$, is the average of m independent samples from $R^\top R$, i.e., $S^\top S = \frac{1}{m} \sum_{i \in [m]} R_i^\top R_i$, for i.i.d. $R_1, R_2, \dots, R_m \sim R$, and therefore,

$$(I_k \times S)^\top (I_k \times S) = \frac{1}{m} \sum_{i \in [m]} (I_k \times R_i)^\top (I_k \times R_i).$$

Recall that $(D_1, D_2) \in \mathcal{E}$ occurs with probability at least $1 - \epsilon\delta/(d\mu_F)$, therefore we have the following for the conditional expectation $\mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid (D_1, D_2) \in \mathcal{E} \right]$,

$$\mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid (D_1, D_2) \in \mathcal{E} \right] \preceq \frac{\mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \right]}{\Pr[(D_1, D_2) \in \mathcal{E}]} \preceq \frac{U^\top U}{1 - \epsilon\delta/(d\mu_F)}.$$

And also by Cauchy-Schwarz we have,

$$\begin{aligned}
&\mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid (D_1, D_2) \in \mathcal{E} \right] \\
&\succeq \mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \right] - \mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid (D_1, D_2) \in \bar{\mathcal{E}} \right] \cdot \Pr[\bar{\mathcal{E}}] \\
&\succeq U^\top U - d \|U\|_F^2 \Pr[\bar{\mathcal{E}}] \cdot I_n \\
&\succeq U^\top U - d \|U\|_F^2 \cdot \epsilon\delta/(d\mu_F) \cdot I_n \\
&\succeq U^\top U - (\epsilon/2) \cdot I_n.
\end{aligned}$$

These two bounds together imply that,

$$\left\| \mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid (D_1, D_2) \in \mathcal{E} \right] - U^\top U \right\|_{op} \leq \epsilon/2.$$

Now note that the random variables $R_i^\top R_i$ are independent conditioned on $(D_1, D_2) \in \mathcal{E}$. Hence, using the upper bounds $L \leq 16\mu_F \cdot \log^2(nd\mu_F k/\epsilon\delta)$ and $M \leq 16\mu_F \mu_2 \cdot \log^2(nd\mu_F k/\epsilon\delta)$, which hold

when $(D_1, D_2) \in \mathcal{E}$, we have the following by Lemma 27, (here we drop the subscript from I_k for ease of notation)

$$\begin{aligned}
& \Pr_{P, D_1, D_2} \left[\left\| U^\top (I \times S)^\top (I \times S) U - U^\top U \right\|_{op} \geq \epsilon \right] \\
& \leq \Pr_P \left[\left\| U^\top (I \times S)^\top (I \times S) U - \mathbb{E} \left[U^\top (I \times R)^\top (I \times R) U \mid (D_1, D_2) \in \mathcal{E} \right] \right\|_{op} \geq \epsilon/2 \mid (D_1, D_2) \in \mathcal{E} \right] \\
& \quad + \Pr_{D_1, D_2} [\bar{\mathcal{E}}] \\
& \leq 8n \cdot \exp \left(-\frac{m\epsilon^2/2}{M + 2\epsilon L/3} \right) + \delta/2 \\
& \leq \delta,
\end{aligned}$$

where the last inequality follows by setting $m = \Omega \left(\log(n/\delta) \log^2(ndk/\epsilon\delta) \cdot \mu_F \mu_2 / \epsilon^2 \right)$. This shows that $I_k \times S$ satisfies the $(\mu_F, \mu_2, \epsilon, \delta, n)$ -spectral property. \square

5.4 Spectral property of Identity \times OSNAP

In this section, we show that tensoring identity operator with OSNAP sketch (Definition 16) results in a transform which satisfies the spectral property (Definition 20) with nearly optimal target dimension as well as nearly optimal application time. This sketch is particularly efficient for sketching sparse vectors. We use a slightly different sketch than the original OSNAP to simplify the analysis, defined as follows.

Definition 21 (OSNAP transform). For every sparsity parameter s , target dimension m , and positive integer d , the OSNAP transform with sparsity parameter s is defined as,

$$S_{r,j} = \sqrt{\frac{1}{s}} \cdot \delta_{r,j} \cdot \sigma_{r,j},$$

for all $r \in [m]$ and all $j \in [d]$, where $\sigma_{r,j} \in \{-1, +1\}$ are independent and uniform Rademacher random variables and $\delta_{r,j}$ are independent Bernoulli random variables satisfying, $\mathbb{E}[\delta_{r,i}] = s/m$ for all $r \in [m]$ and all $i \in [d]$.

Lemma 34. *Suppose $\epsilon, \delta, \mu_2, \mu_F > 0$ and n is a positive integer. If $S \in \mathbb{R}^{m \times d}$ is a OSNAP sketch with sparsity parameter s , then the sketch $I_k \times S$ satisfies the $(\mu_F, \mu_2, \epsilon, \delta, n)$ -spectral property, provided that $s = \Omega \left(\log^2(ndk/\epsilon\delta) \log(n/\delta) \cdot \frac{\mu_2^2}{\epsilon^2} \right)$ and $m = \Omega \left((\mu_F \mu_2 / \epsilon^2) \cdot \log^2(ndk/\epsilon\delta) \right)$.*

Proof. Fix a matrix $U \in \mathbb{R}^{kd \times n}$ with $\|U\|_F^2 \leq \mu_F$ and $\|U\|_{op}^2 \leq \mu_2$. Partition U by rows into $d \times n$ sub-matrices U_1, U_2, \dots, U_k such that $U^T = \begin{bmatrix} U_1^\top & U_2^\top & \cdots & U_k^\top \end{bmatrix}$. Note that

$$U^\top (I_k \times S)^\top (I_k \times S) U = (U_1)^\top S^\top S U_1 + \cdots + (U_k)^\top S^\top S U_k.$$

The proof first considers the simpler case of an OSNAP sketch of rank 1 and then applies the matrix Bernstein bound. Let R denote a rank one OSNAP sketch. R is a $1 \times d$ matrix defined as follows,

$$R_i = \sqrt{\frac{m}{s}} \cdot \delta_i \sigma_i, \tag{28}$$

where σ_i for all $i \in [d]$ are independent Rademacher random variables and also, δ_i for all $i \in [d]$ are independent Bernoulli random variables for which the probability of being one is equal to $\frac{s}{m}$.

We proceed to upper bound the operator norm of $U^\top(I_k \times R)^\top(I_k \times R)U$. First, define the set

$$\mathcal{E} := \left\{ R : (RU_j)^\top RU_j \preceq C \left(\frac{m}{s} \log^2\left(\frac{ndk\mu_F}{\epsilon\delta}\right) \cdot U_j^\top U_j + \log\left(\frac{ndk\mu_F}{\epsilon\delta}\right) \|U_j\|_F^2 \cdot I_n \right) \text{ for all } j = 1, \dots, k \right\},$$

where $C > 0$ is a large enough constant. We show that,

$$\Pr[R \in \mathcal{E}] \geq 1 - \epsilon\delta/(dm\mu_F),$$

where the probability is over the random choices of $\{\sigma_i\}_{i \in [d]}$ and $\{\delta_i\}_{i \in [d]}$. To show this we first prove the following claim,

Claim 35. *For every matrix $Z \in \mathbb{R}^{d \times n}$, if we let R be defined as in (28), then,*

$$\Pr \left[Z^\top R^\top R Z \preceq C \left(\frac{m}{s} \cdot \log^2(n/\delta) Z^\top Z + \log(n/\delta) \|Z\|_F^2 I_n \right) \right] \geq 1 - \delta.$$

Proof. The proof is by Matrix Bernstein inequality, Lemma 26. For any matrix Z let $A = Z(Z^\top Z + \mu I_n)^{-1/2}$, where $\mu = \frac{s}{m} \frac{1}{\log(n/\delta)} \|Z\|_F^2$. We can write $RA = \sqrt{\frac{m}{s}} \sum_{i \in [d]} \delta_i \sigma_i A_i$, where A_i is the i th row of A . Note that $\mathbb{E}[\delta_i \sigma_i A_i] = 0$ and $\|\delta_i \sigma_i A_i\|_2 \leq \|A_i\|_2 \leq \|A\|_{op}$. Also note that

$$\sum_{i \in [d]} \mathbb{E}[(\delta_i \sigma_i A_i)(\delta_i \sigma_i A_i)^*] = \sum_{i \in [d]} \frac{s}{m} \|A_i\|_2^2 = \frac{s}{m} \|A\|_F^2$$

and,

$$\sum_{i \in [d]} \mathbb{E}[(\delta_i \sigma_i A_i)^*(\delta_i \sigma_i A_i)] = \sum_{i \in [d]} \frac{s}{m} A_i^* A_i = \frac{s}{m} A^\top A.$$

Therefore,

$$\max \left\{ \left\| \sum_{i \in [d]} \mathbb{E}[(\delta_i \sigma_i A_i)(\delta_i \sigma_i A_i)^*] \right\|_{op}, \left\| \sum_{i \in [d]} \mathbb{E}[(\delta_i \sigma_i A_i)^*(\delta_i \sigma_i A_i)] \right\|_{op} \right\} \leq \frac{s}{m} \|A\|_F^2.$$

By Lemma 26,

$$\Pr \left[\left\| \sum_{i \in [d]} \delta_i \sigma_i A_i \right\|_{op} \geq t \right] \leq (n+1) \cdot \exp \left(\frac{-t^2/2}{\frac{s}{m} \|A\|_F^2 + \|A\|_{op} t/3} \right),$$

hence if $t = C'/2 \cdot \left(\sqrt{\frac{s}{m} \log(n/\delta)} \|A\|_F + \log(n/\delta) \|A\|_{op} \right)$, then $\Pr \left[\left\| \sum_{i \in [d]} \delta_i \sigma_i A_i \right\|_{op} \geq t \right] \leq \delta$. By plugging $\|RA\|_2^2 = \frac{m}{s} \cdot \left\| \sum_{i \in [d]} \delta_i \sigma_i A_i \right\|_2^2$ into the above we get the following,

$$\Pr \left[\|RA\|_{op}^2 \leq C'^2/2 \left(\frac{m}{s} \cdot \log^2(n/\delta) \|A\|_{op}^2 + \log(n/\delta) \|A\|_F^2 \right) \right] \geq 1 - \delta.$$

Now note that for the choice of $A = Z(Z^\top Z + \mu I_n)^{-1/2}$, we have $\|A\|_{op}^2 \leq \frac{\|Z^\top Z\|_{op}}{\|Z^\top Z\|_{op}^2 + \mu} \leq 1$ and also $\|A\|_F^2 = \sum_i \frac{\lambda_i(Z^\top Z)}{\lambda_i(Z^\top Z) + \mu} \leq \frac{\sum_i \lambda_i(Z^\top Z)}{\mu} = \frac{m}{s} \log(n/\delta)$. By plugging these into the above we get that,

$$\Pr \left[\left\| RZ(Z^\top Z + \mu I_n)^{-1/2} \right\|_{op}^2 \leq C'^2 \frac{m}{s} \cdot \log^2(n/\delta) \right] \geq 1 - \delta.$$

Hence,

$$(Z^\top Z + \mu I_n)^{-1/2} Z^\top R^\top R Z (Z^\top Z + \mu I_n)^{-1/2} \preceq C \frac{m}{s} \cdot \log^2(n/\delta) I_n,$$

with probability $1 - \delta$, where $C = C'^2$. Multiplying both sides of the above from left and right by the positive definite matrix $(Z^\top Z + \mu I_n)^{1/2}$ gives (recall that $\mu = \frac{s}{m} \cdot \frac{\|Z\|_F^2}{\log(n/\delta)}$),

$$Z^\top R^\top R Z \preceq C \left(\frac{m}{s} \cdot \log^2(n/\delta) Z^\top Z + \log(n/\delta) \|Z\|_F^2 I_n \right).$$

□

By applying Claim 35 with failure probability of $\epsilon\delta/(dk\mu_F)$ on each of U_j 's and then applying a union bound, we get the following,

$$\Pr[R \in \mathcal{E}] \geq 1 - \epsilon\delta/(dm\mu_F).$$

From now on, we fix $R \in \mathcal{E}$ and proceed having conditioned on this event.

Upper bounding $\left\| U^\top (I_k \times R)^\top (I_k \times R) U \right\|_{op}$. From the fact that we have conditioned on $R \in \mathcal{E}$, note that,

$$\begin{aligned} L \equiv \left\| U^\top (I_k \times R)^\top (I_k \times R) U \right\|_{op} &= \left\| (U_1)^\top R^\top R U_1 + \dots + (U_k)^\top R^\top R U_k \right\|_{op} \\ &\leq \left\| \sum_{i \in [k]} C \left(\frac{m}{s} \cdot \log^2(ndk\mu_F/\epsilon\delta) \cdot U_j^\top U_j + \log(ndk\mu_F/\epsilon\delta) \|U_j\|_F^2 \cdot I_n \right) \right\|_{op} \\ &= \left\| C \left(\frac{m}{s} \cdot \log^2(ndk\mu_F/\epsilon\delta) \cdot U^\top U + \log(ndk\mu_F/\epsilon\delta) \|U\|_F^2 \cdot I_n \right) \right\|_{op} \\ &\leq C \left(\frac{m}{s} \cdot \log^2(ndk\mu_F/\epsilon\delta) \cdot \|U\|_{op}^2 + \log(ndk\mu_F/\epsilon\delta) \|U\|_F^2 \right) \\ &\leq C \left(\frac{m}{s} \mu_2 \cdot \log^2(ndk\mu_F/\epsilon\delta) + \mu_F \cdot \log(ndk\mu_F/\epsilon\delta) \right). \end{aligned}$$

Upper bounding $\left\| \mathbb{E} \left[\left(U^\top (I_k \times R)^\top (I_k \times R) U \right)^2 \right] \right\|_{op}$. From the condition $R \in \mathcal{E}$, it follows that

$$\begin{aligned} &\mathbb{E} \left[\left(U^\top (I_k \times R)^\top (I_k \times R) U \right)^2 \right] \\ &\preceq \mathbb{E} \left[C \left(\frac{m}{s} \cdot \log^2(ndk\mu_F/\epsilon\delta) \cdot U^\top U + \log(ndk\mu_F/\epsilon\delta) \|U\|_F^2 \cdot I_n \right) \left(U^\top (I_k \times R)^\top (I_k \times R) U \right) \right] \\ &\preceq C \left(\frac{m}{s} \cdot \log^2(ndk\mu_F/\epsilon\delta) \cdot U^\top U + \log(ndk\mu_F/\epsilon\delta) \|U\|_F^2 \cdot I_n \right) \mathbb{E} \left[\left(U^\top (I_k \times R)^\top (I_k \times R) U \right) \right] \\ &\preceq C \left(\frac{m}{s} \cdot \log^2(ndk\mu_F/\epsilon\delta) \cdot U^\top U + \log(ndk\mu_F/\epsilon\delta) \|U\|_F^2 \cdot I_n \right) \cdot \frac{U^\top U}{1 - \epsilon\delta/(dm\mu_F)} \end{aligned}$$

where the last line follows from the fact that the random variable $U^\top (I_k \times R)^\top (I_k \times R) U$ is positive semidefinite and the conditional expectation can be upper bounded by its unconditional expectation as follows,

$$\mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid R \in \mathcal{E} \right] \preceq \frac{\mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \right]}{\Pr[R \in \mathcal{E}]}.$$

Therefore we can bound the operator norm of the above as follows,

$$\begin{aligned}
M &\equiv \left\| \mathbb{E} \left[\left(U^\top (I_k \times R)^\top (I_k \times R) U \right)^2 \right] \right\|_{op} \\
&\leq 2 \left\| C \left(\frac{m}{s} \cdot \log^2(ndk\mu_F/\epsilon\delta) \cdot (U^\top U)^2 + \log(ndk\mu_F/\epsilon\delta) \|U\|_F^2 \cdot U^\top U \right) \right\|_{op} \\
&\leq 2C \left(\frac{m}{s} \cdot \log^2(ndk\mu_F/\epsilon\delta) \cdot \|U^\top U\|_{op}^2 + \log(ndk\mu_F/\epsilon\delta) \|U\|_F^2 \cdot \|U^\top U\|_{op} \right) \\
&= 2C \left(\frac{m}{s} \cdot \log^2(ndk\mu_F/\epsilon\delta) \cdot \mu_2^2 + \log(ndk\mu_F/\epsilon\delta) \mu_F \mu_2 \right).
\end{aligned}$$

Combining one-dimensional OSNAP transforms. To conclude, we note that the Gram matrix of an OSNAP sketch, $S^\top S \in \mathbb{R}^{d \times d}$, is the average of m independent samples from $R^\top R$ with R defined as in (28) – i.e., $S^\top S = \frac{1}{m} \sum_{i \in [m]} R_i^\top R_i$ for i.i.d. $R_1, R_2, \dots, R_m \sim R$, and therefore,

$$(I_k \times S)^\top (I_k \times S) = \frac{1}{m} \sum_{i \in [m]} (I_k \times R_i)^\top (I_k \times R_i).$$

Note that by a union bound $R_i \in \mathcal{E}$ simultaneously for all $i \in [m]$ with probability at least $1 - \epsilon\delta/(d\mu_F)$. Now note that the random variables $R_i^\top R_i$ are independent conditioned on $R_i \in \mathcal{E}$ for all $i \in [m]$. Also note that the conditional expectation $\mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid R \in \mathcal{E} \right]$ satisfies the following,

$$\begin{aligned}
&\mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid R \in \mathcal{E} \right] \\
&\geq \mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \right] - \mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid R \in \bar{\mathcal{E}} \right] \cdot \Pr[\bar{\mathcal{E}}] \\
&\geq U^\top U - d \|U\|_F^2 \Pr[\bar{\mathcal{E}}] \cdot I_n \\
&\geq U^\top U - d \|U\|_F^2 \cdot \epsilon\delta/(d\mu_F) \cdot I_n \\
&\geq U^\top U - d \|U\|_F^2 \cdot \epsilon/2 \cdot I_n.
\end{aligned}$$

We also have,

$$\mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid R \in \mathcal{E} \right] \preceq \frac{\mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \right]}{\Pr[R \in \mathcal{E}]} \preceq \frac{U^\top U}{1 - \epsilon\delta/(d\mu_F)}.$$

These two bounds together imply that,

$$\left\| \mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid R \in \mathcal{E} \right] - U^\top U \right\|_{op} \leq \epsilon/2.$$

Now, using the upper bounds $L \leq C \left(\frac{m}{s} \mu_2 \cdot \log^2(ndk\mu_F/\epsilon\delta) + \mu_F \cdot \log(ndk\mu_F/\delta) \right)$ and $M \leq 2C \left(\frac{m}{s} \cdot \log^2(ndk\mu_F/\delta) \cdot \mu_2^2 + \log(ndk\mu_F/\delta) \mu_F \mu_2 \right)$, which hold when $R \in \mathcal{E}$, we have that by Lemma 27,

$$\begin{aligned}
&\Pr \left[\left\| U^\top (I_k \times S)^\top (I_k \times S) U - U^\top U \right\|_{op} \geq \epsilon \right] \\
&\leq \Pr \left[\left\| U^\top (I_k \times S)^\top (I_k \times S) U - \mathbb{E} \left[U^\top (I_k \times R)^\top (I_k \times R) U \mid R \in \mathcal{E} \right] \right\|_{op} \geq \epsilon/2 \mid \mathcal{E} \right] + \Pr[\bar{\mathcal{E}}] \\
&\leq 8n \cdot \exp \left(-\frac{m\epsilon^2/8}{M + \epsilon L/3} \right) + \delta/2 \leq \delta,
\end{aligned}$$

where the last inequality follows by setting $s = \Omega\left(\log^2(ndk\mu_F/\epsilon\delta)\log(nd/\delta)\cdot\frac{\mu_2^2}{\epsilon^2}\right)$ and $m = \Omega\left(\mu_F\mu_2/\epsilon^2\cdot\log^2(ndk\mu_F/\epsilon\delta)\right)$. This shows that $I_k \times S$ satisfies the $(\mu_F, \mu_2, \epsilon, \delta, n)$ -spectral property. \square

5.5 High Probability OSE with linear dependence on s_λ

We are ready to prove Theorem 3. We prove that if we instantiate Π^p from Definition 11 with $T_{\text{base}} : \text{OSNAP}$ and $S_{\text{base}} : \text{TensorSRHT}$, it satisfies the statement of Theorem 3.

Theorem 3. *For every positive integers p, d, n , every $\epsilon, s_\lambda > 0$, there exists a distribution on linear sketches $\Pi^p \in \mathbb{R}^{m \times d^p}$ which is an $(\epsilon, 1/\text{poly}(n), s_\lambda, d^p, n)$ -oblivious subspace embedding as in Definition 2, provided that the integer m satisfies $m = \tilde{\Omega}(p^4 s_\lambda / \epsilon^2)$.*

Moreover, for any $X \in \mathbb{R}^{d \times n}$, if $A \in \mathbb{R}^{d^p \times n}$ is the matrix whose columns are obtained by a p -fold self-tensoring of each column of X then the matrix $\Pi^p A$ can be computed using Algorithm 1 in time $\tilde{O}(pnm + p^5 \epsilon^{-2} \text{nnz}(X))$.

Proof. Let $\delta = \frac{1}{\text{poly}(n)}$ denote the failure probability. Let $m \approx p^4 \log_2^3(\frac{nd}{\epsilon\delta}) \cdot \frac{s_\lambda}{\epsilon^2}$ and $s \approx \frac{p^4}{\epsilon^2} \cdot \log_2^3(\frac{nd}{\epsilon\delta})$ be integers. Let $\Pi^p \in \mathbb{R}^{m \times m^p}$ be the sketch defined in Definition 11, where $S_{\text{base}} \in \mathbb{R}^{m \times m^2}$ is a TensorSRHT sketch and $T_{\text{base}} \in \mathbb{R}^{m \times d}$ is an OSNAP sketch with sparsity parameter s .

Let $q = 2^{\lceil \log_2(p) \rceil}$. By Lemma 10, it is sufficient to show that Π^q is a $(\epsilon, \delta, s_\lambda, d^q, n)$ -Oblivious Subspace Embedding. Consider arbitrary $A \in \mathbb{R}^{d^q \times n}$ and $\lambda > 0$. Let us denote the statistical dimension of A by $s_\lambda = s_\lambda(A^\top A)$. Let $U = A(A^\top A + \lambda I_n)^{-1/2}$. Therefore, $\|U\|_2 \leq 1$ and $\|U\|_F^2 = s_\lambda$. Since $q < 2p$, by Lemma 34, the transform $I_{d^{q-1}} \times T_{\text{base}}$, satisfies $(2s_\lambda + 2, 2, O(\epsilon/q), O(\delta/n^2q), n)$ -spectral property. Moreover, by Lemma 33, the transform $I_{m^{q-2}} \times S_{\text{base}}$ satisfies $(5s_\lambda + 9, 9, O(\epsilon/q), O(\delta/n^2q^2), n)$ -spectral property. Therefore, by Lemma 32, the sketch Π^q satisfies $(s_\lambda + 1, 1, \epsilon, \delta, n)$ -spectral property, hence,

$$\Pr \left[\left\| (\Pi^q U)^\top \Pi^q U - U^\top U \right\|_{\text{op}} \leq \epsilon \right] \geq 1 - \delta.$$

Since $U^\top U = (A^\top A + \lambda I_n)^{-1/2} A^\top A (A^\top A + \lambda I_n)^{-1/2}$ and $\Pi^q U = \Pi^p A (A^\top A + \lambda I_n)^{-1/2}$ we have the following,

$$\Pr \left[(1 - \epsilon)(A^\top A + \lambda I_n) \preceq (\Pi^p A)^\top \Pi^p A + \lambda I_n \preceq (1 + \epsilon)(A^\top A + \lambda I_n) \right] \geq 1 - \delta.$$

Runtime: By Lemma 8, for any S_{base} and T_{base} , if A is the matrix whose columns are obtained by p -fold self-tensoring of each column of some $X \in \mathbb{R}^{d \times n}$ then the sketched matrix $\Pi^p A$ can be computed using Algorithm 1. When S_{base} is TensorSRHT and T_{base} is OSNAP, the runtime of Algorithm 1 for a fixed vector $w \in \mathbb{R}^d$ is as follows; Computing Y_j^0 's for each j in lines 3 and 4 of algorithm requires applying an OSNAP sketch on $w \in \mathbb{R}^d$ which on expectation takes time $O(s \cdot \text{nnz}(w))$. Therefore computing all Y_j^0 's takes time $O(qs \cdot \text{nnz}(w))$.

Computing each of Y_j^l 's in line 7 of algorithm amounts to applying a TensorSRHT of input dimension m^2 and target dimension of m on $Y_{2j-1}^{l-1} \otimes Y_{2j}^{l-1}$. This takes time $O(m \log m)$. Therefore computing all the Y_j^l 's takes time $O(q \cdot m \log m)$. Note that $q \leq 2p$ hence the total time of running Algorithm 1 on a vector w is $O(p \cdot m \log_2 m + ps \cdot \text{nnz}(w))$. Therefore, sketching n columns of a matrix $X \in \mathbb{R}^{d \times n}$ takes time $O(p(nm \log_2 m + s \cdot \text{nnz}(X)))$. \square

6 Oblivious Subspace Embedding for the Gaussian Kernel

In this section we show how to sketch the Gaussian kernel matrix by polynomial expansion and then applying our proposed sketch for the polynomial kernels.

Data-points with bounded ℓ_2 radius: Suppose that we are given a dataset of points $x_1, \dots, x_n \in \mathbb{R}^d$ such that for all $i \in [n]$, $\|x_i\|_2^2 \leq r$ for some positive value r . Consider the Gaussian kernel matrix $G \in \mathbb{R}^{n \times n}$ defined as $G_{i,j} = e^{-\|x_i - x_j\|_2^2/2}$ for all $i, j \in [n]$. We are interested in sketching the data-points matrix X using a sketch $S_g : \mathbb{R}^d \rightarrow \mathbb{R}^m$ such that the following holds with probability $1 - \delta$,

$$(1 - \epsilon)(G + \lambda I_n) \preceq (S_g(X))^\top S_g(X) + \lambda I_n \preceq (1 + \epsilon)(G + \lambda I_n).$$

Theorem 5. *For every $r > 0$, every positive integers n, d , and every $X \in \mathbb{R}^{d \times n}$ such that $\|x_i\|_2 \leq r$ for all $i \in [n]$, where x_i is the i^{th} column of X , suppose $G \in \mathbb{R}^{n \times n}$ is the Gaussian kernel matrix – i.e., $G_{j,k} = e^{-\|x_j - x_k\|_2^2/2}$ for all $j, k \in [n]$. There exists an algorithm which computes $S_g(X) \in \mathbb{R}^{m \times n}$ in time $\tilde{O}(q^6 \epsilon^{-2} n s_\lambda + q^6 \epsilon^{-2} \text{nnz}(X))$ such that for every $\epsilon, \lambda > 0$,*

$$\Pr_{S_g} \left[(1 - \epsilon)(G + \lambda I_n) \preceq (S_g(X))^\top S_g(X) + \lambda I_n \preceq (1 + \epsilon)(G + \lambda I_n) \right] \geq 1 - 1/\text{poly}(n),$$

where $m = \tilde{\Theta}(q^5 s_\lambda / \epsilon^2)$ and $q = \Theta(r^2 + \log(n/\epsilon\lambda))$ and s_λ is λ -statistical dimension of G as in Definition 1.

Proof. Let $\delta = \frac{1}{\text{poly}(n)}$ denote the failure probability. Note that $G_{i,j} = e^{-\|x_i\|_2^2/2} \cdot e^{x_i^\top x_j} \cdot e^{-\|x_j\|_2^2/2}$ for every $i, j \in [n]$. Let D be a $n \times n$ diagonal matrix with i th diagonal entry $e^{-\|x_i\|_2^2/2}$ and let $K \in \mathbb{R}^{n \times n}$ be defined as $K_{i,j} = e^{x_i^\top x_j}$ (note that $DKD = G$). Note that K is a positive definite kernel matrix. The Taylor series expansion for kernel K is as follows,

$$K = \sum_{l=0}^{\infty} \frac{(X^{\otimes l})^\top X^{\otimes l}}{l!}.$$

Therefore G can be written as the following series,

$$G = \sum_{l=0}^{\infty} \frac{(X^{\otimes l} D)^\top X^{\otimes l} D}{l!}.$$

Note that each of the terms $(X^{\otimes l} D)^\top X^{\otimes l} D = D(X^{\otimes l})^\top X^{\otimes l} D$ are positive definite kernel matrices. The statistical dimension of kernel $(X^{\otimes l} D)^\top X^{\otimes l} D$ for every $l \geq 0$ is upper bounded by the statistical dimension of kernel G through the following claim.

Claim 36. *For every $\mu \geq 0$ and every integer l ,*

$$s_\mu \left((X^{\otimes l} D)^\top X^{\otimes l} D \right) \leq s_\mu(G).$$

Proof. From the Taylor expansion $G = \sum_{l=0}^{\infty} \frac{(X^{\otimes l} D)^\top X^{\otimes l} D}{l!}$ along with the fact that the polynomial kernel of any degree is positive definite, we have that $(X^{\otimes l} D)^\top X^{\otimes l} D \preceq G$. Now, by Courant-Fischer's min-max theorem we have that,

$$\lambda_j((X^{\otimes l} D)^\top X^{\otimes l} D) = \max_{U \in \mathbb{R}^{(j-1) \times n}} \min_{\substack{\alpha \neq 0 \\ U\alpha=0}} \frac{\alpha^\top (X^{\otimes l} D)^\top X^{\otimes l} D \alpha}{\|\alpha\|_2^2}.$$

Let U^* be the maximizer of the expression above. Then we have,

$$\begin{aligned}
\lambda_j(G) &= \max_{U \in \mathbb{R}^{(j-1) \times n}} \min_{\substack{\alpha \neq 0 \\ U\alpha=0}} \frac{\alpha^\top G \alpha}{\|\alpha\|_2^2} \\
&\geq \min_{\substack{\alpha \neq 0 \\ U^* \alpha = 0}} \frac{\alpha^\top G \alpha}{\|\alpha\|_2^2} \\
&\geq \min_{\substack{\alpha \neq 0 \\ U^* \alpha = 0}} \frac{\alpha^\top (X^{\otimes l} D)^\top X^{\otimes l} D \alpha}{\|\alpha\|_2^2} \\
&= \lambda_j((X^{\otimes l} D)^\top X^{\otimes l} D).
\end{aligned}$$

for all j . Therefore, the claim follows from the definition of statistical dimension,

$$s_\mu(G) = \sum_{j=1}^n \frac{\lambda_j(G)}{\lambda_j(G) + \mu} \geq \sum_{j=1}^n \frac{\lambda_j((X^{\otimes l} D)^\top X^{\otimes l} D)}{\lambda_j((X^{\otimes l} D)^\top X^{\otimes l} D) + \mu} = s_\mu((X^{\otimes l} D)^\top X^{\otimes l} D).$$

□

If we let $P = \sum_{l=0}^q \frac{(X^{\otimes l})^\top X^{\otimes l}}{l!}$, where $q = C \cdot (r^2 + \log(\frac{n}{\epsilon\lambda}))$ for some constant C , then by the triangle inequality we have

$$\begin{aligned}
\|K - P\|_{op} &\leq \sum_{l>q} \left\| \frac{(X^{\otimes l})^\top X^{\otimes l}}{l!} \right\|_{op} \\
&\leq \sum_{l>q} \left\| \frac{(X^{\otimes l})^\top X^{\otimes l}}{l!} \right\|_F \\
&\leq \sum_{l>q} \frac{n \cdot r^{2l}}{l!} \\
&\leq \epsilon\lambda/2.
\end{aligned}$$

P is a positive definite kernel matrix. Also note that all the eigenvalues of the diagonal matrix D are bounded by 1. Hence, in order to get a subspace embedding it is sufficient to satisfy the following with probability $1 - \delta$,

$$(1 - \epsilon/2)(DPD + \lambda I_n) \preceq (S_g(X))^\top S_g(X) + \lambda I_n \preceq (1 + \epsilon/2)(DPD + \lambda I_n).$$

Let the sketch $\Pi^l \in \mathbb{R}^{m_l \times d^l}$ be the sketch from Theorem 3 therefore by Claim 36 we get the following guarantee on Π^l :

$$(1 - \frac{\epsilon}{9})((X^{\otimes l} D)^\top X^{\otimes l} D + \lambda I_n) \preceq (\Pi^l X^{\otimes l} D)^\top \Pi^l X^{\otimes l} D + \lambda I_n \preceq (1 + \frac{\epsilon}{9})((X^{\otimes l} D)^\top X^{\otimes l} D + \lambda I_n), \quad (29)$$

with probability $1 - \frac{\delta}{q+1}$ as long as $m_l = \Omega\left(l^4 \log^3(nd/\delta) \cdot s_\lambda/\epsilon^2\right)$ and moreover $\Pi^l X^{\otimes l} D$ can be computed using $O\left(n \cdot l \cdot m_l \log_2 m_l + \frac{l^5}{\epsilon^2} \cdot \log^3(nd/\delta) \cdot \text{nnz}(X)\right)$ runtime where s_λ is the λ -statistical dimension of G .

We let S_P be the sketch of size $m \times (\sum_{l=0}^q d^l)$ which sketches the kernel P . The sketch S_P is defined as

$$S_P = \frac{1}{\sqrt{0!}}\Pi^0 \oplus \frac{1}{\sqrt{1!}}\Pi^1 \oplus \frac{1}{\sqrt{2!}}\Pi^2 \cdots \frac{1}{\sqrt{q!}}\Pi^q.$$

Let Z be the matrix of size $(\sum_{l=0}^q d^l) \times n$ whose i^{th} column is

$$z_i = x_i^{\otimes 0} \oplus x_i^{\otimes 1} \oplus x_i^{\otimes 2} \cdots x_i^{\otimes q},$$

where x_i is the i^{th} column of X . Therefore the following holds for $(S_P Z)^\top S_P Z$,

$$(S_P Z)^\top S_P Z = \sum_{l=0}^q \frac{(\Pi^l X^{\otimes l})^\top \Pi^l X^{\otimes l}}{l!},$$

and hence,

$$(S_P Z D)^\top S_P Z D = \sum_{l=0}^q \frac{(\Pi^l X^{\otimes l} D)^\top \Pi^l X^{\otimes l} D}{l!}.$$

Therefore by combining the terms of (29) for all $0 \leq l \leq q$, using a union bound we get that with probability $1 - \delta$, the following holds,

$$(1 - \epsilon/2)(DPD + \lambda I_n) \preceq (S_P Z D)^\top S_P Z D + \lambda I_n \preceq (1 + \epsilon/2)(DPD + \lambda I_n).$$

Now we define $S_g(x)$ which is a non-linear transformation on the input x defined as

$$S_g(x) = e^{-\|x\|_2^2/2} \left(\frac{1}{\sqrt{0!}} \cdot \Pi^0(x^{\otimes 0}) \oplus \frac{1}{\sqrt{1!}} \cdot \Pi^1(x^{\otimes 1}) \oplus \frac{1}{\sqrt{2!}} \cdot \Pi^2(x^{\otimes 2}) \cdots \frac{1}{\sqrt{q!}} \cdot \Pi^q(x^{\otimes q}) \right).$$

We have that $S_g(X) = S_P Z D$, therefore with probability $1 - \delta$, the following holds,

$$(1 - \epsilon)(G + \lambda I_n) \preceq (S_g(X))^\top S_g(X) + \lambda I_n \preceq (1 + \epsilon)(G + \lambda I_n).$$

Note that the target dimension of S_g is $m = m_0 + m_1 + \cdots + m_q \approx q^5 \log^3(nd/\delta) s_\lambda / \epsilon^2$. Also, by Theorem 3, time to compute $S_g(X)$ is $O\left(\frac{nq^6}{\epsilon^2} \cdot \log^4(nd/\delta) \cdot s_\lambda + \frac{q^6}{\epsilon^2} \cdot \log^3(nd/\delta) \cdot \text{nnz}(X)\right)$. \square

Acknowledgements

Michael Kapralov is supported by ERC Starting Grant SUBLINEAR. Thomas D. Ahle, Jakob B. T. Knudsen, and Rasmus Pagh are supported by Villum Foundation grant 16582 to Basic Algorithms Research Copenhagen (BARC). David Woodruff is supported in part by Office of Naval Research (ONR) grant N00014-18-1-2562. Part of this work was done while Michael Kapralov, Rasmus Pagh, and David Woodruff were visiting the Simons Institute for the Theory of Computing.

References

- [AC06] Nir Ailon and Bernard Chazelle. Approximate nearest neighbors and the fast johnson-lindenstrauss transform. In *Proceedings of the 38th Annual ACM Symposium on Theory of Computing, Seattle, WA, USA, May 21-23, 2006*, pages 557–563, 2006.
- [Ach03] Dimitris Achlioptas. Database-friendly random projections: Johnson-lindenstrauss with binary coins. *J. Comput. Syst. Sci.*, 66(4):671–687, June 2003.

- [ACW17a] Haim Avron, Kenneth L. Clarkson, and David P. Woodruff. Faster kernel ridge regression using sketching and preconditioning. *SIAM J. Matrix Analysis Applications*, 38(4):1116–1138, 2017.
- [ACW17b] Haim Avron, Kenneth L. Clarkson, and David P. Woodruff. Sharper bounds for regularized data fitting. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM 2017, August 16-18, 2017, Berkeley, CA, USA*, pages 27:1–27:22, 2017.
- [AK19] Thomas D Ahle and Jakob BT Knudsen. Almost optimal tensor sketch. *arXiv preprint arXiv:1909.01821*, 2019.
- [AKM⁺17] Haim Avron, Michael Kapralov, Cameron Musco, Christopher Musco, Ameya Velingker, and Amir Zandieh. Random fourier features for kernel ridge regression: Approximation bounds and statistical guarantees. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 253–262, 2017.
- [AKM⁺18a] Haim Avron, Michael Kapralov, Cameron Musco, Christopher Musco, Ameya Velingker, and Amir Zandieh. Random fourier features for kernel ridge regression: Approximation bounds and statistical guarantees. *CoRR*, abs/1804.09893, 2018.
- [AKM⁺18b] Haim Avron, Michael Kapralov, Cameron Musco, Christopher Musco, Ameya Velingker, and Amir Zandieh. A universal sampling method for reconstructing signals with simple fourier transforms. *arXiv preprint arXiv:1812.08723*, 2018.
- [AM15] Ahmed El Alaoui and Michael W. Mahoney. Fast randomized kernel ridge regression with statistical guarantees. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 775–783, 2015.
- [ANW14] Haim Avron, Huy Nguyen, and David Woodruff. Subspace embeddings for the polynomial kernel. In *Advances in neural information processing systems*, pages 2258–2266, 2014.
- [BCL⁺10] Vladimir Braverman, Kai-Min Chung, Zhenming Liu, Michael Mitzenmacher, and Rafail Ostrovsky. AMS without 4-wise independence on product domains. In *27th International Symposium on Theoretical Aspects of Computer Science, STACS 2010, March 4-6, 2010, Nancy, France*, pages 119–130, 2010.
- [BLM13] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- [CCFC02] Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams. In *International Colloquium on Automata, Languages, and Programming*, pages 693–703. Springer, 2002.
- [Cha02] Moses S Charikar. Similarity estimation techniques from rounding algorithms. In *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*, pages 380–388. ACM, 2002.

- [CJN18] Michael B. Cohen, T. S. Jayram, and Jelani Nelson. Simple analyses of the sparse johnson-lindenstrauss transform. In *1st Symposium on Simplicity in Algorithms, SOSA 2018, January 7-10, 2018, New Orleans, LA, USA*, pages 15:1–15:9, 2018.
- [CKS11] Andrew Cotter, Joseph Keshet, and Nathan Srebro. Explicit approximations of the gaussian kernel. *arXiv preprint arXiv:1109.4603*, 2011.
- [CNW16a] Michael B. Cohen, Jelani Nelson, and David P. Woodruff. Optimal approximate matrix product in terms of stable rank. In *43rd International Colloquium on Automata, Languages, and Programming, ICALP 2016, July 11-15, 2016, Rome, Italy*, pages 11:1–11:14, 2016.
- [CNW16b] Michael B. Cohen, Jelani Nelson, and David P. Woodruff. Optimal approximate matrix product in terms of stable rank. In *43rd International Colloquium on Automata, Languages, and Programming, ICALP 2016, July 11-15, 2016, Rome, Italy*, pages 11:1–11:14, 2016.
- [Coh16] Michael B. Cohen. Nearly tight oblivious subspace embeddings by trace inequalities. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2016, Arlington, VA, USA, January 10-12, 2016*, pages 278–287, 2016.
- [CW09] Kenneth L. Clarkson and David P. Woodruff. Numerical linear algebra in the streaming model. In *Proceedings of the 41st Annual ACM Symposium on Theory of Computing, STOC 2009, Bethesda, MD, USA, May 31 - June 2, 2009*, pages 205–214, 2009.
- [CW13] Kenneth L. Clarkson and David P. Woodruff. Low rank approximation and regression in input sparsity time. In *Symposium on Theory of Computing Conference, STOC’13, Palo Alto, CA, USA, June 1-4, 2013*, pages 81–90, 2013.
- [CW17] Kenneth L. Clarkson and David P. Woodruff. Low-rank approximation and regression in input sparsity time. *Journal of the ACM (JACM)*, 63(6):54, 2017.
- [DKS10] Anirban Dasgupta, Ravi Kumar, and Tamás Sarlós. A sparse johnson: Lindenstrauss transform. In *Proceedings of the 42nd ACM Symposium on Theory of Computing, STOC 2010, Cambridge, Massachusetts, USA, 5-8 June 2010*, pages 341–350, 2010.
- [DIPG12] Victor De la Pena and Evarist Giné. *Decoupling: from dependence to independence*. Springer Science & Business Media, 2012.
- [DMM06] Petros Drineas, Michael W. Mahoney, and S. Muthukrishnan. Sampling algorithms for l_2 regression and applications. In *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2006, Miami, Florida, USA, January 22-26, 2006*, pages 1127–1136, 2006.
- [DMMS11] Petros Drineas, Michael W. Mahoney, S. Muthukrishnan, and Tamás Sarlós. Faster least squares approximation. *Numerische Mathematik*, 117(2):219–249, 2011.
- [DMMW12] Petros Drineas, Malik Magdon-Ismail, Michael W. Mahoney, and David P. Woodruff. Fast approximation of matrix coherence and statistical leverage. *Journal of Machine Learning Research*, 13:3475–3506, 2012.

- [Hit93] Paweł Hitczenko. Domination inequality for martingale transforms of a rademacher sequence. *Israel Journal of Mathematics*, 84(1-2):161–178, 1993.
- [Hit94] Paweł Hitczenko. On a domination of sums of random variables by sums of conditionally independent ones. *The Annals of Probability*, pages 453–468, 1994.
- [HM07] Uffe Haagerup and Magdalena Musat. On the best constants in noncommutative khintchine-type inequalities. *Journal of Functional Analysis*, 250(2):588–624, 2007.
- [IM08] Piotr Indyk and Andrew McGregor. Declaring independence via the sketching of sketches. In *Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 737–745. Society for Industrial and Applied Mathematics, 2008.
- [JLS86] William B. Johnson, Joram Lindenstrauss, and Gideon Schechtman. Extensions of lipschitz maps into banach spaces. *Israel Journal of Mathematics*, 54(2):129–138, Jun 1986.
- [KN14] Daniel M. Kane and Jelani Nelson. Sparser johnson-lindenstrauss transforms. *J. ACM*, 61(1):4:1–4:23, 2014.
- [KPV⁺19] Michael Kapralov, Rasmus Pagh, Ameya Velingker, David Woodruff, and Amir Zandieh. Oblivious sketching of high-degree polynomial kernels. *arXiv preprint arXiv:1909.01410*, 2019.
- [KVW14] Ravi Kannan, Santosh Vempala, and David Woodruff. Principal component analysis and higher correlations for distributed data. In *Conference on Learning Theory*, pages 1040–1057, 2014.
- [Lat97] Rafał Łatała. Estimation of moments of sums of independent real random variables. *The Annals of Probability*, 25(3):1502–1513, 1997.
- [Lat06] Rafał Łatała. Estimates of moments and tails of gaussian chaoses. *The Annals of Probability*, 34(6):2315–2331, 2006.
- [LDFU13] Yichao Lu, Paramveer S. Dhillon, Dean P. Foster, and Lyle H. Ungar. Faster ridge regression via the subsampled randomized hadamard transform. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.*, pages 369–377, 2013.
- [LSS14] Quoc Viet Le, Tamás Sarlós, and Alexander Johannes Smola. Fastfood: Approximate kernel expansions in loglinear time. *CoRR*, abs/1408.3060, 2014.
- [MM17] Cameron Musco and Christopher Musco. Recursive sampling for the nystrom method. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 3836–3848, 2017.
- [NDT15] Nam H Nguyen, Petros Drineas, and Trac D Tran. Tensor sparsification via a bound on the spectral norm of random tensors. *Information and Inference: A Journal of the IMA*, 4(3):195–229, 2015.

- [NN13] 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.
- [Pag13] Rasmus Pagh. Compressed matrix multiplication. *TOCT*, 5(3):9:1–9:17, 2013.
- [PP13] Ninh Pham and Rasmus Pagh. Fast and scalable polynomial kernels via explicit feature maps. In *The 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2013, Chicago, IL, USA, August 11-14, 2013*, pages 239–247, 2013.
- [PT12] Mihai Patrascu and Mikkel Thorup. The power of simple tabulation hashing. *J. ACM*, 59(3):14:1–14:50, 2012.
- [PW15] Mert Pilanci and Martin J. Wainwright. Randomized sketches of convex programs with sharp guarantees. *IEEE Trans. Information Theory*, 61(9):5096–5115, 2015.
- [RR07] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 3-6, 2007*, pages 1177–1184, 2007.
- [RR08] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2008.
- [Sar06] Tamás Sarlós. Improved approximation algorithms for large matrices via random projections. In *47th Annual IEEE Symposium on Foundations of Computer Science (FOCS 2006), 21-24 October 2006, Berkeley, California, USA, Proceedings*, pages 143–152, 2006.
- [Tro11] Joel A. Tropp. Improved analysis of the subsampled randomized hadamard transform. *Advances in Adaptive Data Analysis*, 3(1-2):115–126, 2011.
- [Tro15] Joel A. Tropp. An introduction to matrix concentration inequalities. *Foundations and Trends in Machine Learning*, 8(1-2):1–230, 2015.
- [Val15] Gregory Valiant. Finding correlations in subquadratic time, with applications to learning parities and the closest pair problem. *Journal of the ACM (JACM)*, 62(2):13, 2015.
- [Woo14] David P. Woodruff. Sketching as a tool for numerical linear algebra. *Foundations and Trends in Theoretical Computer Science*, 10(1-2):1–157, 2014.

A Direct Lower and Upper Bounds

We introduce the following notation. We say $f(x) \lesssim g(x)$ if for some universal constant C we have $f(x) \leq Cg(x)$ for all $x \in \mathcal{R}$ and \cdot . Note this is slightly different from the usual $f(x) = O(g(x))$ in that it is uniform in x rather than asymptotic. We similarly say $f(x) \gtrsim g(x)$ if $g(x) \lesssim f(x)$ and $f(x) \sim g(x)$ if both $f(x) \lesssim g(x)$ and $f(x) \gtrsim g(x)$.

We will also make heavy use of the L^p norm notation for random variables in \mathcal{R} , that is for $p \geq 1$ we write $\|X\|_{L^p} = (E|X|^p)^{1/p}$. A very useful result for computing the L^p -norm of a sum of random variables is the following:

Lemma 37 (Latala's inequality, [Lat97]). *If $p \geq 2$ and X, X_1, \dots, X_n are iid. mean 0 random variables, then we have*

$$\left\| \sum_{i=1}^n X_i \right\|_{L^p} \sim \sup \left\{ \frac{p}{s} \left(\frac{n}{p} \right)^{1/s} \|X\|_{L^s} \mid \max \left\{ 2, \frac{p}{n} \right\} \leq s \leq p \right\}. \quad (30)$$

The following simple corollary will be used for both upper and lower bounds:

Corollary 38. *Let $p \geq 2, C > 0$ and $\alpha \geq 1$. Let $(X_i)_{i \in [n]}$ be iid. mean 0 random variables such that $\|X_i\|_{L^p} \sim (Cp)^\alpha$, then $\|\sum_i X_i\|_{L^p} \sim C^\alpha \max\{2^\alpha \sqrt{pn}, (n/p)^{1/p} p^\alpha\}$.*

Proof. We will show that the expression in eq. (30) is maximized either by minimizing or maximizing s . Hence we need to check that $\frac{p}{s} \left(\frac{n}{p} \right)^{1/s} s^\alpha$ it has no other optimums in the valid range. For this, we note that $\frac{d}{ds} \frac{p}{s} \left(\frac{n}{p} \right)^{1/s} s^\alpha = \frac{-p}{s^{3-\alpha}} \left(\frac{n}{p} \right)^{1/s} \left((1-\alpha)s + \log \frac{n}{p} \right)$. Given $\alpha \geq 1$ the derivative is non-decreasing in s , which gives the lemma. \square

For the lower bound we will also use the following result by Hitzzenko, which provides an improvement on Khintchine for Rademacher random variables.

Lemma 39 (Sharp bound on Rademacher sums [Hit93]). *Let $\sigma \in \{-1, 1\}^n$ be a random Rademacher sequence and let $a \in \mathbb{R}^n$ be an arbitrary real vector with sorted entries $|a_1| \geq |a_2| \geq \dots \geq |a_n|$, then*

$$\|\langle a, \sigma \rangle\|_{L^p} \sim \sum_{i \leq p} a_i + \sqrt{p} \left(\sum_{i > p} a_i^2 \right)^{1/2} \quad (31)$$

Finally the lower bound will use the Paley-Zygmund inequality (also known as the one-sided Chebyshev inequality):

Lemma 40 (Paley-Zygmund). *Let $X \geq 0$ be a real random variable with finite variance, and let $\theta \in [0, 1]$, then*

$$\Pr[X \geq \theta \mathbb{E}[X]] \geq (1 - \theta)^2 \frac{\mathbb{E}[X]^2}{\mathbb{E}[X^2]}. \quad (32)$$

A classical strategy when using Paley-Zygmund is to prove $\mathbb{E}[X] \geq 2\varepsilon$ for some $\varepsilon > 0$, and then take $\theta = 1/2$ to give $\Pr[X \geq \varepsilon] \geq \mathbb{E}[X]^2 / (4\mathbb{E}[X^2])$.

A.1 Lower Bound for Sub-Gaussians

The following lower bound considers the sketching matrix consisting of the direct composition of matrices with Rademacher entries. Note however that the assumptions on Rademachers are only used to show that the p -norm of a single row with a vector is $\sim \sqrt{p}$. For this reason the same lower bound hold if the Rademacher entries are substituted for, say Gaussians.

Theorem 41 (Lower bound). *For some constants $C_1, C_2, B > 0$, let $d, m, c \geq 1$ be integers, let $\varepsilon \in [0, 1]$ and $\delta \in [0, 1/16]$. Further assume that $d \geq \log 1/\delta \geq c/B$. Then the following holds.*

Let $M^{(1)}, \dots, M^{(c)} \in \mathcal{R}^{m \times d}$ be matrices with all independent Rademacher entries and let $M = \frac{1}{\sqrt{m}} M^{(1)} \bullet \dots \bullet M^{(c)}$. Then there exists some unit vector $y \in \mathcal{R}^{dc}$ such that if

$$m < C_1 \max \left\{ 3^c \varepsilon^{-2} \frac{\log 1/\delta}{c}, \varepsilon^{-1} \left(\frac{C_2 \log 1/\delta}{c} \right)^c \right\} \quad \text{then} \quad \Pr \left[\left| \|My\|_2^2 - 1 \right| > \varepsilon \right] > \delta. \quad (33)$$

Proof. Let $y = [1, \dots, 1]^T / \sqrt{d} \in \mathcal{R}^d$ and let $x = y^{\otimes c}$. We have

$$\|My\|_2^2 - 1 = \frac{1}{m} \left\| M^{(1)}x \circ \dots \circ M^{(c)}x \right\|_2^2 - 1 = \frac{1}{m} \sum_{j \in [m]} \left(\prod_{i \in [c]} Z_{i,j}^2 \right) - 1 \quad (34)$$

where each $Z_{i,j} = \sum_{k \in [d]} M_{j,k}^{(i)} / \sqrt{d}$ are independent averages of d independent Rademacher random variables. By Lemma 39 we have $\|Z_{i,j}\|_{L^p} \sim \min\{\sqrt{p}, \sqrt{d}\}$ which is \sqrt{p} by the assumption $d \geq \log 1/\delta$ as long as $p \leq \log 1/\delta$. By the expanding $Z_{i,j}^4$ into monomials and linearity of expectation we get $\|Z_{i,j}\|_4 = \frac{1}{\sqrt{d}}(d + 3d(d-1))^{1/4} = (3 - 2/d)^{1/4}$.

Now define $X_j = \prod_{i \in [c]} Z_{i,j}^2 - 1$, then $EX_j = 0$ and $\|X_j\|_{L^p} \geq \left\| \prod_{i \in [c]} Z_{i,j}^2 \right\|_{L^p} - 1 = \|Z_{i,j}\|_{L^{2p}}^{2c} - 1 \geq K^c p^c$ for some K , assuming $p \geq 2$. In particular, $\|X_j\|_{L^2} \geq \|Z_{i,j}\|_{L^4}^{2c} - 1 = (3 - 2/d)^{c/2} - 1 \sim 3^{c/2}$ by the assumption $d \geq c \geq 1$.

We have $\| \|My\|_2^2 - 1 \|_{L^p} = \frac{1}{m} \left\| \sum_{j \in [m]} X_j \right\|_{L^p}$ is a sum of iid. random variables, so we can use Corollary 38 to show

$$K_3 \max \left\{ \sqrt{3^c p/m}, (m/p)^{1/p} K_1^c p^c/m \right\} \lesssim \left\| \|My\|_2^2 - 1 \right\|_{L^p} \quad (35)$$

$$\lesssim K_4 \max \left\{ \sqrt{3^c p/m}, (m/p)^{1/p} K_2^c p^c/m \right\} \quad (36)$$

for some universal constants $K_1, K_2, K_3, K_4 > 0$.

Assume now that $m < \max \left\{ AK_3^2 3^c \varepsilon^{-2} \frac{\log 1/\delta}{c}, \frac{K_3}{4} \varepsilon^{-1} \left(4AK_1 \frac{\log 1/\delta}{c} \right)^c \right\}$ as in the theorem. We take $p = 4A \frac{\log 1/\delta}{c}$ for some constant A to be determined. We want to show $\| \|My\|_2^2 - 1 \|_{L^p} \geq 2\varepsilon$. For this we split into two cases depending on which term of $m < \max\{(1), (2)\}$ dominates. If (1) \geq (2) we pick the first lower bound in eq. (36) and get $\| \|My\|_2^2 - 1 \|_{L^p} \geq K_3 \sqrt{3^c p/m} \geq K_3 \sqrt{\frac{4\varepsilon^2}{K_3^2}} = 2\varepsilon$. Otherwise, if (2) \geq (1), we pick the other lower bound and also get:

$$\left\| \|My\|_2^2 - 1 \right\|_{L^p} \geq K_3 (m/p)^{1/p} \frac{K_1^c p^c}{m} \geq \frac{K_3}{2} \frac{K_1^c \left(4A \frac{\log 1/\delta}{c} \right)^c}{\frac{K_3}{4} \varepsilon^{-1} \left(4AK_1 \frac{\log 1/\delta}{c} \right)^c} = 2\varepsilon, \quad (37)$$

where we used $(m/p)^{1/p} \geq e^{-1/(em)} \geq 1/2$ for $m \geq 1$. Plugging into Paley-Zygmund (Lemma 40) we have

$$\Pr \left[\left| \|My\|_2^2 - 1 \right| \geq \varepsilon \right] \geq \Pr \left[\left| \|My\|_2^2 - 1 \right|^p \geq \left\| \|My\|_2^2 - 1 \right\|_{L^p}^p 2^{-p} \right] \quad (38)$$

$$\geq \frac{1}{4} \left(\frac{\left\| \|My\|_2^2 - 1 \right\|_{L^p}^p}{\left\| \|My\|_2^2 - 1 \right\|_{L^{2p}}^p} \right)^{2p}, \quad (39)$$

where we used that $p \geq 1$ so $(1 - 2^{-p})^2 \geq 1/4$.

There are again two cases depending on which term of the upper bound in eq. (36) dominates. If $\sqrt{3^c p/m} \geq (m/p)^{1/p} K_2^c p^c/m$ we have using the first lower bound that $\frac{\left\| \|My\|_2^2 - 1 \right\|_{L^p}}{\left\| \|My\|_2^2 - 1 \right\|_{L^{2p}}} \geq \frac{K_3}{\sqrt{2}K_4}$.

For the alternative case, $(m/p)^{1/p} K_2^c p^c/m \geq \sqrt{3^c p/m}$, we have

$$\frac{\left\| \|My\|_2^2 - 1 \right\|_{L^p}}{\left\| \|My\|_2^2 - 1 \right\|_{L^{2p}}} \geq \frac{K_3}{\sqrt{2}K_4} \frac{(m/p)^{1/p}}{(m/2p)^{1/2p}} \left(\frac{K_1}{2K_2} \right)^c \geq \frac{K_3}{2K_4} \left(\frac{K_1}{2K_2} \right)^c \quad (40)$$

where $\frac{(m/p)^{1/p}}{(m/2p)^{1/2p}} \geq e^{-1/(4em)} \geq 1/\sqrt{2}$ for $m \geq 1$.

Comparing with (39) we see that it suffices to take $A \leq \min\{\frac{1}{\log 2K_4/K_3}, \frac{1}{\log 2K_2/K_1}\}/32$. This choice also ensures that $1 \leq p \leq \log 1/\delta$ as we promised. Note that we may assume in eq. (36) that $K_3 \leq K_4$ and $K_1 \leq K_2$. We then finally have

$$\frac{1}{4} \left(\frac{K_3}{\sqrt{2}K_4} \right)^{2p} \geq \frac{1}{4} \delta^{1/(4c)} \quad \text{and} \quad \frac{1}{4} \left(\frac{K_3}{2K_4} \left(\frac{K_1}{2K_2} \right)^c \right)^{2p} \geq \frac{1}{4} \delta^{1/(4c)+1/4}, \quad (41)$$

which are both $\geq \delta$ for $c \geq 1$ and $\delta < 1/16$. \square

A.2 Upper bound for Sub-Gaussians

Theorem 42 (Upper bound). *Let $\varepsilon, \delta \in [0, 1]$ and let $\gamma > 0$, $1 \leq c \leq \frac{\log 1/\delta}{4\gamma}$ be some constants. Let $T \in \mathbb{R}^{m \times d}$ be a matrix with iid. rows $T_1, \dots, T_m \in \mathcal{R}^d$ such that $\mathbb{E}[(T_1 x)^2] = \|x\|_2^2$ and $\|T_1 x\|_{L^p} \leq \sqrt{ap} \|x\|_2$ for some $a > 0$ and $p \geq 4$. Let $M = T^{(1)} \bullet \dots \bullet T^{(c)}$ where $T^{(1)}, \dots, T^{(c)}$ are independent copies of T . Then M has the JL-moment property, $\| \|Mx\|_2 - \|x\|_2 \|_{L^p} \leq \varepsilon \delta^{1/p}$, given*

$$m \gtrsim (4ae^\gamma)^{2c} \varepsilon^{-2} \frac{\log 1/\delta}{c\gamma} + (4ae^\gamma)^c \varepsilon^{-1} \left(\frac{\log 1/\delta}{c\gamma} \right)^c. \quad (42)$$

Remark 2. In the case of random Rademachers we set $a = \sqrt{3}/4$ to get

$$m = O \left(3^c \varepsilon^{-2} \frac{\log 1/\delta}{c\gamma} e^{2c\gamma} + \varepsilon^{-1} \left(\sqrt{3} \frac{\log 1/\delta}{c\gamma} \right)^c e^{c\gamma} \right).$$

Note that depending on γ this matches either of the terms of the lower bound. Setting $\gamma = \Theta(1/c)$ or $\gamma = \Theta(1)$ we have either

$$m = O \left(3^c \varepsilon^{-2} \log 1/\delta + \varepsilon^{-1} \left(\sqrt{3} \log 1/\delta \right)^c \right) \quad \text{or} \quad m = O \left((3e^2)^c \varepsilon^{-2} \frac{\log 1/\delta}{c} + \varepsilon^{-1} \left(\sqrt{3} e \frac{\log 1/\delta}{c} \right)^c \right).$$

Finally, in the case of constant $c = O(1)$, $\gamma = \Theta(1)$ we simply get

$$m = O \left(\varepsilon^{-2} \log 1/\delta + \varepsilon^{-1} (\log 1/\delta)^c \right).$$

Proof of Theorem 42. Without loss of generalization we may assume $\|x\|_{L^2} = 1$. We notice that $\| \|Mx\|_2^2 - 1 \|_{L^p} \leq \left\| \frac{1}{m} \sum_i (M_i x)^2 - 1 \right\|_{L^p}$ is the mean of iid. random variables. Call these $Z_i = (M_i x)^2 - 1$. Then $EZ_i = 0$ and $\|Z_i\|_{L^p} = \|(M_i x)^2 - 1\|_{L^p} \lesssim \|(M_i x)^2\|_{L^p} = \|M_i x\|_{L^{2p}}^2$ by symmetrization. Now by the assumption $\|T_1 x\|_{L^p} \leq \sqrt{ap} \|x\|_2 = \sqrt{ap}$, and by Lemma 19, we get that $\|M_i x\|_{L^p} = \|T_i^{(1)} \otimes \dots \otimes T_i^{(c)} x\|_{L^p} \leq (ap)^{c/2}$, and so $\|Z_i\|_{L^p} \leq (2ap)^c$ for all $i \in [m]$.

We now use Corollary 38 which implies

$$\left\| \frac{1}{m} \sum_i Z_i \right\|_{L^p} \lesssim (4a)^c \sqrt{p/m} + m^{1/p} (2ap)^c / m \lesssim (4a)^c \sqrt{p/m} + (4ap)^c / m. \quad (43)$$

The second inequality comes from the following consideration: If the second term of (43) dominates, then $(4a)^c \sqrt{p/m} \leq m^{1/p} (2ap)^c / m$ which implies $m^{1/p} \leq (p/2)^{\frac{2c-1}{p-2}} \leq 2^c$ for $p \geq 4$.

All that remains is to decide on p . We take $p = \frac{\log 1/\delta}{c\gamma}$ which is ≥ 4 by assumption, and $m = \max\{(4ae^\gamma)^{2c}p\varepsilon^{-2}, (4ae^\gamma)^c p^c \varepsilon^{-1}\}$. Then

$$\left\| \frac{1}{m} \sum_i Z_i \right\|_{L^p}^p \lesssim (4a)^{cp} \max\{\varepsilon^p (4ae^\gamma)^{-cp}, \varepsilon^p (4ae^\gamma)^{-cp}\} \quad (44)$$

$$= e^{-c\gamma p} \varepsilon^p \quad (45)$$

$$= \delta \varepsilon^p, \quad (46)$$

which is exactly the JL moment property. \square

A.3 Lower Bound for TensorSketch

For every integer d, q , the TensorSketch of degree q , $M : \mathbb{R}^{d^q} \rightarrow \mathbb{R}^m$ is defined as,

$$M(x^{\otimes q}) = \mathcal{F}^{-1}((\mathcal{F}C_1x) \circ (\mathcal{F}C_2x) \circ \cdots \circ (\mathcal{F}C_qx)), \quad (47)$$

for every $x \in \mathbb{R}^d$ where $C_1, \dots, C_q \in \mathbb{R}^{m \times d}$ are independent instances of CountSketch and $\mathcal{F} \in \mathbb{C}^{m \times m}$ is the Discrete Fourier Transform matrix with proper normalization which satisfies the convolution theorem, also note that, \circ denotes entry-wise (Hadamard) product of vectors of the same size.

Lemma 43. *For every integer d, q , let $M : \mathbb{R}^{d^q} \rightarrow \mathbb{R}^m$ be the TensorSketch of degree $q \leq d$, see (47). For the all ones vector $x = \{1\}^d$,*

$$\text{Var} \left[\|Mx^{\otimes q}\|_2^2 \right] \geq \left(\frac{3^q}{2m^2} - 1 \right) \|x^{\otimes q}\|_2^4.$$

Proof. Note that since \mathcal{F} is normalized such that it satisfies the convolution theorem, \mathcal{F}^{-1} is indeed a unitary matrix times $1/\sqrt{m}$, $\|Mx^{\otimes q}\|_2^2 = \frac{1}{m} \|(\mathcal{F}C_1x) \circ (\mathcal{F}C_2x) \circ \cdots \circ (\mathcal{F}C_qx)\|_2^2$. Consider the first entry of the vector $(\mathcal{F}C_1x) \circ (\mathcal{F}C_2x) \circ \cdots \circ (\mathcal{F}C_qx)$. Because the first row of \mathcal{F} is all ones $\{1\}^m$, the first element of the mentioned vector for the choice of $x = \{1\}^d$ is $\prod_{i=1}^q \left(\sum_{j \in [d]} \sigma^i(j) \right) = \prod_{i=1}^q \left(\sum_{j \in [d]} \sigma^i(j) \right)$, where $\sigma^i : [d] \rightarrow \{-1, +1\}$ are fully independent random hash functions used by the CountSketch C_i for all $i \in [q]$. Let us denote by V the following positive random variable,

$$V = \prod_{i=1}^q \left(\sum_{j \in [d]} \sigma^i(j) \right)^2.$$

Note that $\|Mx^{\otimes q}\|_2^2 \geq \frac{V}{m}$, hence $\mathbb{E}[\|Mx^{\otimes q}\|_2^4] \geq \frac{\mathbb{E}[V^2]}{m^2}$. Also note that $\mathbb{E}[V^2] = \prod_{i=1}^q \mathbb{E} \left[\left(\sum_{j \in [d]} \sigma^i(j) \right)^4 \right]$ because σ^i 's are independent. We can write

$$\mathbb{E} \left[\left(\sum_{j \in [d]} \sigma^i(j) \right)^4 \right] = 3d^2 - 2d = 3 \left(1 - \frac{1}{6d} \right) \|x\|_2^4,$$

hence if $d \geq q$,

$$\mathbb{E}[V^2] \geq (1/2) \cdot 3^q \cdot \|x^{\otimes q}\|_2^4,$$

Therefore $\mathbb{E}[\|Mx^{\otimes q}\|_2^4] \geq \frac{\mathbb{E}[V^2]}{m^2} \geq \frac{3^q}{2m^2} \|x^{\otimes q}\|_2^4$. It is also true that $\mathbb{E}[\|Mx^{\otimes q}\|_2^2] = \|x^{\otimes q}\|_2^2$ [ANW14]. \square

Lemma 44. For every integer d, q every $\varepsilon > 0$, every $0 < \delta \leq \frac{1}{2 \cdot 12^q}$, let $M : \mathbb{R}^{d^q} \rightarrow \mathbb{R}^m$ be the TensorSketch of degree q , see (47). If $m < 3^{q/2}$ then for the all ones vector $x = \{1\}^d$ we have,

$$\Pr \left[\left| \|Mx^{\otimes q}\|_2^2 - \|x^{\otimes q}\|_2^2 \right| > 1/2 \cdot \|x^{\otimes q}\|_2^2 \right] > \delta.$$

Proof. Note that since \mathcal{F} is normalized such that it satisfies the convolution theorem, \mathcal{F}^{-1} is indeed a unitary matrix times $1/\sqrt{m}$, $\|Mx^{\otimes q}\|_2^2 = \frac{1}{m} \|(\mathcal{F}C_1x) \circ (\mathcal{F}C_2x) \circ \dots \circ (\mathcal{F}C_qx)\|_2^2$. Consider the first entry of the vector $(\mathcal{F}C_1x) \circ (\mathcal{F}C_2x) \circ \dots \circ (\mathcal{F}C_qx)$. Because the first row of \mathcal{F} is all ones $\{1\}^m$, the first element of the mentioned vector for the choice of $x = \{1\}^d$ is $\prod_{i=1}^q \left(\sum_{j \in [d]} \sigma^i(j) \right) = \prod_{i=1}^q \left(\sum_{j \in [d]} \sigma^i(j) \right)$, where $\sigma^i : [d] \rightarrow \{-1, +1\}$ are fully independent random hash functions used by the CountSketch C_i for all $i \in [q]$. Let us denote by V the following positive random variable,

$$V = \prod_{i=1}^q \left(\sum_{j \in [d]} \sigma^i(j) \right)^2.$$

Note that $\|Mx^{\otimes q}\|_2^2 \geq \frac{V}{m}$. Note that $\mathbb{E}[V^t] = \prod_{i=1}^q \mathbb{E} \left[\left(\sum_{j \in [d]} \sigma^i(j) \right)^{2t} \right]$ for every t because σ^i 's are independent. Note that for $t = 2$ we have,

$$\mathbb{E} \left[\left(\sum_{j \in [d]} \sigma^i(j) \right)^4 \right] = 3d^2 - 2d \geq 3 \left(1 - \frac{1}{6d}\right) \|x\|_2^4,$$

hence if $d \geq q$,

$$\mathbb{E}[V^2] \geq (3^q/2) \cdot \|x^{\otimes q}\|_2^4.$$

Now consider $t = 4$. By Khintchine's inequality, Lemma 17, we have,

$$\mathbb{E} \left[\left(\sum_{j \in [d]} \sigma^i(j) \right)^8 \right] \leq 105 \cdot \|x\|_2^8,$$

hence,

$$\mathbb{E}[V^4] \leq 105^q \cdot \|x^{\otimes q}\|_2^8.$$

Therefore by Paley Zygmund we have the following,

$$\begin{aligned} \Pr \left[\|Mx^{\otimes q}\|_2^2 \geq \frac{3^{\frac{q}{2}}}{2m} \cdot \|x^{\otimes q}\|_2^2 \right] &\geq \Pr \left[V \geq 3^{\frac{q}{2}}/2 \cdot \|x^{\otimes q}\|_2^2 \right] \\ &= \Pr \left[V^2 \geq 3^q/4 \cdot \|x^{\otimes q}\|_2^4 \right] \\ &\geq \Pr \left[V^2 \geq 1/4 \cdot \mathbb{E}[V^2] \right] \\ &\geq 1/2 \cdot \frac{\mathbb{E}[V^2]^2}{\mathbb{E}[V^4]} \\ &\geq \frac{9^q}{2 \cdot 105^q} \\ &> \frac{1}{2 \cdot 12^q} \geq \delta. \end{aligned}$$

