
Sketching for Kronecker Product Regression and P-splines

<p>Huaian Diao hadiao@nenu.edu.cn Northeast Normal University</p>	<p>Zhao Song zhaos@utexas.edu Harvard U & UT-Austin</p>	<p>Wen Sun wensun@cs.cmu.edu Carnegie Mellon University</p>	<p>David P. Woodruff dwoodruf@cs.cmu.edu Carnegie Mellon University</p>
--	--	--	--

Abstract

TENSORSKETCH is an oblivious linear sketch introduced in (Pagh, 2013) and later used in (Pham and Pagh, 2013) in the context of SVMs for polynomial kernels. It was shown in (Avron et al., 2014) that TENSORSKETCH provides a subspace embedding, and therefore can be used for canonical correlation analysis, low rank approximation, and principal component regression for the polynomial kernel. We take TENSORSKETCH outside of the context of polynomial kernels, and show its utility in applications in which the underlying design matrix is a Kronecker product of smaller matrices. This allows us to solve Kronecker product regression and non-negative Kronecker product regression, as well as regularized spline regression. Our main technical result is then in extending TENSORSKETCH to other norms. That is, TENSORSKETCH only provides input sparsity time for Kronecker product regression with respect to the 2-norm. We show how to solve Kronecker product regression with respect to the 1-norm in time sublinear in the time required for computing the Kronecker product, as well as for more general p -norms.

1 INTRODUCTION

In the overconstrained least squares regression problem, we are given an $n \times d$ matrix A called the “design matrix”, $n \gg d$, and an $n \times 1$ vector b , and the goal is to find an x which minimizes $\|Ax - b\|_2$. There are many variants to this problem, such as regularized versions of the problem in which one seeks an x so as to

minimize $\|Ax - b\|_2^2 + \|Lx\|_2^2$ for a matrix L , or regression problems which seek to minimize more robust loss functions, such as ℓ_1 -regression $\|Ax - b\|_1$.

In the era of big data, large scale matrix computations have attracted considerable interest. To obtain reasonable computational and time complexities for large scale matrix computations, a number of randomized approximation algorithms have been developed. For example, in (Woodruff, 2014), it was shown how to output a vector $x' \in \mathbb{R}^d$ for which $\|Ax' - b\|_2 \leq (1 + \epsilon) \min_{x \in \mathbb{R}^d} \|Ax - b\|_2$ in $\text{nnz}(A) + \text{poly}(d/\epsilon)$ time, where $\text{nnz}(A)$ denotes the number of non-zero entries of matrix A . We refer the reader to the recent surveys (Kannan and Vempala, 2009; Mahoney, 2011; Woodruff, 2014) for a detailed treatment of this topic.

In this work we focus on regression problems for which the design matrix is a *Kronecker product matrix*, that is, it has the form $A \otimes B$ for $A \in \mathbb{R}^{n_1 \times d_1}$ and $B \in \mathbb{R}^{n_2 \times d_2}$. Also, $b \in \mathbb{R}^{n_1 n_2}$, and one seeks to solve the problem $\min_{x \in \mathbb{R}^{d_1 d_2}} \|(A \otimes B)x - b\|_2$ in the standard setting, which can also be generalized to regularized and robust variants. One can also ask the question when the design matrix is a Kronecker product of more than two matrices.

Kronecker product matrices have many applications in statistics, linear system theory, signal processing, photogrammetry, multivariate data fitting, etc.; see (Golub and Van Loan, 2013; Van Loan and Pitsianis, 1993; Van Loan, 1992). The linear least squares problem involving the Kronecker product arises in many applications, such as structured linear total least norm on blind deconvolution problems (Oh and Yun, 2005), constrained linear least squares problems with a Kronecker product structure, the bivariate problem of surface fitting and multidimensional density smoothing (Eilers and Marx, 2006).

One way to solve Kronecker product regression is to form the matrix $C = A \otimes B$ explicitly, where $A \in \mathbb{R}^{n_1 \times d_1}$, $B \in \mathbb{R}^{n_2 \times d_2}$, and then apply a randomized algorithm to C . However, this takes at least $\text{nnz}(A) \cdot \text{nnz}(B)$ time and space. It is natural

Proceedings of the 21st International Conference on Artificial Intelligence and Statistics (AISTATS) 2018, Lanzarote, Spain. JMLR: W&CP volume 7X. Copyright 2018 by the author(s).

to ask if it is possible to solve the Kronecker product regression problem in time faster than computing $A \otimes B$. This is in fact the case, as Fausett and Fulton (Fausett and Fulton, 1994) show that one can solve Kronecker product regression in $O(n_1 d_1^2 + n_2 d_2^2)$ time; indeed, the solution vector $x = \text{vec}((B^\perp)^\top D^{-1} A^\perp)$, where $D = \text{vec}(b)$ and $\text{vec}(E)$ for a matrix E denotes the operation of stacking the columns of E into a single long vector. While such a computation does not involve computing $A \otimes B$, it is more expensive than what one would like.

A natural question is if one can approximately solve Kronecker product regression in $\text{nnz}(A) + \text{nnz}(B) + \text{poly}(d/\epsilon)$ time, which, up to the $\text{poly}(d/\epsilon)$ term, would match the description size of the input. Another natural question is Kronecker product regression with regularization. Such regression problems arise frequently in the context of splines (Eilers and Marx, 2006). Finally, what about Kronecker product regression with other, more robust, norms such as the ℓ_1 -norm?

1.1 Our Contributions

We first observe that the random linear map TENSORSKETCH, introduced in the context of problems for the polynomial kernel by Pagh (2013), Pham and Pagh (2013), is exactly suited for this task. Namely, in (Avron et al., 2014) it was shown that for a d -dimensional subspace C of \mathbb{R}^n , represented as an $n \times d$ matrix, there is a distribution on linear maps S with $O(d^2/\epsilon^2)$ rows such that with constant probability, simultaneously for all $x \in \mathbb{R}^d$, $\|SCx\|_2 = (1 \pm \epsilon)\|Cx\|_2$. That is, S provides an *Oblivious Subspace Embedding (OSE)* for the column span of C . Further, it is known that if b is an n -dimensional vector, then one has that $\|S[C, b]x\|_2 = (1 \pm \epsilon)\|[C, b]x\|_2$ for all $x \in \mathbb{R}^{d+1}$. Consequently, to solve the regression problem $\min_{x \in \mathbb{R}^d} \|Cx - b\|_2$, one can instead solve the much smaller problem $\min_{x \in \mathbb{R}^d} \|SCx - Sb\|_2$, and the solution x' to the latter problem will satisfy $\|Cx' - b\|_2 \leq (1 + \epsilon) \min_{x \in \mathbb{R}^d} \|Cx - b\|_2$.

Importantly, if $n = n_1 \cdot n_2$ and there is a basis for the column span of C of the form $A_1 \otimes A_1, A_2 \otimes A_2, \dots, A_d \otimes A_d$, then given A_1, \dots, A_d , it holds that SC can be computed in $\text{nnz}(A)$ time, where A is the $n_1 \times d$ matrix whose i -th column is A_i . Further, given a vector b with $\text{nnz}(b)$ non-zero entries, one can compute Sb in $\text{nnz}(b)$ time. Thus, one obtains a $(1 + \epsilon)$ -approximate solution to the regression $\min_{x \in \mathbb{R}^d} \|Cx - b\|_2$ in $\text{nnz}(A) + \text{nnz}(b) + \text{poly}(d/\epsilon)$ time.

While not immediately useful for our problem, we show that via simple modifications, the claim about TENSORSKETCH above can be generalized to the case when there is a basis for the column span of C of the form $A_i \otimes B_j$ for arbitrary vectors $A_1, \dots, A_{d_1} \in$

\mathbb{R}^{n_1} and vectors $B_1, \dots, B_{d_2} \in \mathbb{R}^{n_2}$. That is, we observe that in this case SC can be computed in $\text{nnz}(A) + \text{nnz}(B)$ time, where A is the $n_1 \times d_1$ matrix whose i -th column is A_i , and B is the $n_2 \times d_2$ matrix whose i -th column is B_i . In this case $C = A \otimes B$, which is exactly the case of Kronecker product regression. Using the above connection to regression, we obtain an algorithm for Kronecker product regression in $\text{nnz}(A) + \text{nnz}(B) + \text{nnz}(b) + \text{poly}(d_1 d_2/\epsilon)$ time. Using the fact that $\|SCx - Sb\|_2 = (1 \pm \epsilon)\|Cx - b\|_2$ for all x , we have in particular that this holds for all non-negative x , and so also obtain the same reduction in problem size for non-negative Kronecker product regression, which occurs often in image and text data; see e.g., (Chen and Plemmons, 2010).

The above observation allows us to extend many existing variants of least squares regression to the case when C is a Kronecker product of two matrices. For example, the results in (Avron et al., 2016) for the ridge regression problem $\min_{x \in \mathbb{R}^d} \|Ax - b\|_2^2 + \lambda \|x\|_2^2$ immediately hold when C is a Kronecker product matrix, since the conditions needed for the oblivious embedding in (Avron et al., 2016) come down to a subspace embedding and an approximate matrix product condition, both of which are known to hold for TENSORSKETCH (Avron et al., 2014). More interestingly we are able to extend the results in (Avron et al., 2016) for Kronecker ridge regression to the case when the regularizer is a general matrix, namely, to the problem $\min_{x \in \mathbb{R}^d} \|Ax - b\|_2^2 + \lambda \|Lx\|_2^2$, where L is an arbitrary matrix. Such problems occur in the context of spline regression (Eilers and Marx, 1996, 2006; Eilers et al., 2015). The number of rows of TENSORSKETCH depends on a generalized notion of statistical dimension depending on the generalized singular values γ_i of $[A; L]$, and may be much smaller than d_1 or d_2 . In (Avron et al., 2016), only results for L equal to the identity were obtained.

Finally, our main technical result is to extend our results to least absolute deviation Kronecker product regression $\min_{x \in \mathbb{R}^d} \|Cx - b\|_1$, which, since it involves the 1-norm, is often considered more robust than least squares regression (Rousseeuw and Leroy, 2005) and has been widely used in applications such as computer vision (Zheng et al., 2012). We in fact extend this to general p -norms but focus the discussion on $p = 1$. We give the first algorithms for this problem that are faster than computing $C = A \otimes B$ explicitly, which would take at least $\text{nnz}(A) \cdot \text{nnz}(B) \geq n_1 n_2$ time. Namely, and for simplicity focusing on the case when $n_1 = n_2$, for which the goal is to do better than n_1^2 time, we show how to output an $x \in \mathbb{R}^{d_1 \times d_2}$ for which $\|Cx' - b\|_1 \leq (1 + \epsilon) \min_{x \in \mathbb{R}^{d_1 d_2}} \|Cx - b\|_1$ in time $n^{3/2} \text{poly}(d_1 d_2/\epsilon)$. While this is more expensive than

solving Kronecker product least squares, the 1-norm may lead to significantly more robust solutions. From a technical perspective, TENSORSKETCH when applied to a vector actually destroys the 1-norm of a vector, preserving only its 2-norm, so new ideas are needed. We show how to use multiple TENSORSKETCH matrices to implicitly obtain very crude estimates to the so-called ℓ_1 -leverage scores of C , which can be interpreted as probabilities of sampling rows of C in order to obtain an ℓ_1 -subspace embedding of the column span of C (see, e.g., (Woodruff, 2014)). However, since C has $n_1 n_2$ rows, we cannot even afford to write down the ℓ_1 -leverage scores of C . We show how to implicitly represent such leverage scores and to sample from them without ever explicitly writing them down. Balancing the phases of our algorithm leads to our overall time.

1.2 Notation

We consider Kronecker Product of q 2-d matrices $A_1 \otimes A_2 \otimes \dots \otimes A_q$, where each $A_i \in \mathbb{R}^{n_i \times d_i}$. We denote $\mathcal{A} = A_1 \otimes A_2 \otimes \dots \otimes A_q$, $n = \prod_{i=1}^q n_i$, and $d = \prod_{i=1}^q d_i$. We denote $[n]$ as the set $\{1, 2, 3, \dots, n\}$. The ℓ_p norm for a vector $x \in \mathbb{R}^d$ is defined as $\|x\|_p = (\sum_{i=1}^d |x_i|^p)^{1/p}$, where x_i stands for the i 'th entry of the vector x . For any matrix M , we use $M_{i,*}$ to represent the i 'th row of M and $M_{*,j}$ as the j 'th column of M .

We define a *Well-Conditioned Basis* and *Statistical Dimension* as follow (similar definitions can be found in Clarkson (2005); Dasgupta et al. (2009); Sohler and Woodruff (2011); Meng and Mahoney (2013); Song et al. (2017a,b) and Avron et al. (2016)):

Definition 1.1 (Well-Conditioned Basis). *Let A be an $n \times m$ matrix with rank d , let $p \in [1, \infty)$, and let $\|\cdot\|_q$ be the dual norm of $\|\cdot\|_p$, i.e., $1/p + 1/q = 1$. Then an $n \times d$ matrix U is an (α, β, p) -well-conditioned basis for the column space of A , if the columns of U span the column space of A and (1) $\|U\|_p \leq \alpha$, and (2) for all $z \in \mathbb{R}^d$, $\|z\|_q \leq \beta \|Uz\|_p$.*

Consider the classic Ridge Regression: $\min_x \|Ax - b\|_2^2 + \lambda \|x\|_2^2$. The *Statistical Dimension* is defined as:

Definition 1.2 (Statistical Dimension). *Let A be an $n \times m$ matrix with rank d and singular values $\sigma_i, i \in [d]$. For $\lambda \geq 0$, the statistical dimension $\text{sd}_\lambda(A)$ is defined as the quantity $\text{sd}_\lambda(A) = \sum_{i \in [d]} 1/(1 + \lambda/\sigma_i^2)$.*

2 BACKGROUND:TENSORSKETCH

We briefly introduce TENSORSKETCH (Pagh, 2013; Avron et al., 2014) and how to apply TENSORSKETCH to the Kronecker product of multiple matrices efficiently without explicitly computing the tensor product.¹

¹We refer readers to (Avron et al., 2014) for more details about TENSORSKETCH.

We want to find a oblivious subspace embedding S such that for any $x \in \mathbb{R}^n$, we have $\|SAx\|_2 = (1 \pm \epsilon)\|Ax\|_2$, where the notation $a = (1 \pm \epsilon)b$ stands for $(1 - \epsilon)b \leq a \leq (1 + \epsilon)b$, for any $a, b \in \mathbb{R}$. Consider the (i_1, i_2, \dots, i_q) 'th column of \mathcal{A} ($i_j \in [d_j]$): $A_{1*,i_1} \otimes A_{2*,i_2} \otimes \dots \otimes A_{q*,i_q}$. Assume the sketching target dimension is m . TENSORSKETCH is defined using q 3-wise independent hash functions $h_i : [n_i] \rightarrow [m]$, and q 4-wise independent sign functions $s_i : [n_i] \rightarrow \{-1, 1\}$, $\forall i \in [q]$. Define hash function $H : [n_1] \times [n_2] \times \dots \times [n_q] \rightarrow [m]$ as $H(i_1, i_2, \dots, i_q) = ((\sum_{k=1}^q h_k(i_k)) \bmod m)$, and sign function $S : [n_1] \times [n_2] \times \dots \times [n_q] \rightarrow \{-1, 1\}$ as $S(i_1, i_2, \dots, i_q) = \prod_{k=1}^q s_k(i_k)$. Applying TENSORSKETCH to the Kronecker product of vectors $A_{1*,i_1}, \dots, A_{q*,i_q}$ is equivalent to applying COUNTSKETCH (Charikar et al., 2004) defined with H and S to the vector $(A_{1*,i_1} \otimes \dots \otimes A_{q*,i_q})$. To apply TENSORSKETCH to \mathcal{A} , we just need to apply COUNTSKETCH defined with H and S to the columns of \mathcal{A} one by one.

Applying TENSORSKETCH to the Kronecker product of $(A_{1*,i_1}, \dots, A_{q*,i_q})$ naively would require at least $O(n)$ time. Pagh (2013) shows that one can apply TENSORSKETCH to the Kronecker product of $(A_{1*,i_1}, \dots, A_{q*,i_q})$ without explicitly computing the Kronecker product of these vectors using the Fast Fourier Transformation. Particularly, Pham and Pagh (2013) show that one only needs $O(\sum_{j=1}^q \text{nnz}(A_{j*,i_j}) + qm \log(m))$ time to compute $S(A_{1*,i_1} \otimes \dots \otimes A_{q*,i_q})$ where $A_{i*,j}$ stands for the j 'th column of A_i . As \mathcal{A} has d columns, computing $S\mathcal{A}$ takes $O(d(\sum_{i=1}^q \text{nnz}(A_i) + dqm \log(m)))$ time, which is much smaller than $O(\prod_{i=1}^q n_i)$, which is the least amount of time one needs for explicitly computing \mathcal{A} . In the rest of the paper, we assume that we compute $S\mathcal{A}$ using the above efficient procedure without explicitly computing $A_1 \otimes \dots \otimes A_q$.

3 TENSOR PRODUCT LEAST SQUARES REGRESSION

Consider the tensor product least squares regression problem $\min_x \|Ax - b\|_2$ where $\mathcal{A} \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$. Let $S \in \mathbb{R}^{m \times n}$ be the matrix form of TENSORSKETCH of Section 2. We propose a TENSORSKETCH -type Algorithm 1 for the tensor product regression problem. The following theorem shows that the solution obtained from Alg. 1 is a good approximation of the optimal solution of the original tensor product regression.

Let us define OPT to be the optimal cost of the optimization problem, e.g., $\text{OPT} = \min_x \|Ax - b\|_2$. The following theorem shows that Alg. 1 computes a good approximate solution.

Theorem 3.1. (Tensor regression) Suppose \tilde{x} is the output of Algorithm 1 with TENSORSKETCH $S \in$

Algorithm 1 Tensor product regression

- 1: **procedure** TREGRESSION(A, b, ϵ, δ)
 - 2: $m \leftarrow (d_1 d_2 \cdots d_q + 1)^2 (2 + 3^q) / (\epsilon^2 \delta)$
 - 3: Choose S to be an $m \times (n_1 n_2 \cdots n_q)$ TENSORSKETCH matrix
 - 4: Compute $S(A_1 \otimes A_2 \otimes \cdots \otimes A_q)$ and Sb
 - 5: $\tilde{x} \leftarrow \min_x \|S(A_1 \otimes A_2 \otimes \cdots \otimes A_q)x - Sb\|_2$
 - 6: **return** \tilde{x}
 - 7: **end procedure**
-

$\mathbb{R}^{m \times n}$, where $m = 8(d_1 d_2 \cdots d_q + 1)^2 (2 + 3^q) / (\epsilon^2 \delta)$. Then the following approximation $\|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)\tilde{x} - b\|_2 \leq (1 + \epsilon) \text{OPT}$, holds with probability at least $1 - \delta$.

The proof of Theorem 3.1 can be found in Appendix C.1. Theorem 3.1 shows that we can achieve an ϵ -close solution by solving a much smaller regression problem with a number of samples of order $O(\text{poly}(d/\epsilon))$, which is independent of the large dimension n . Using the technique we introduced in Sec. 2, we can also compute SA without explicitly computing the tensor product.

We can extend Theorem 3.1 to the nonnegative tensor product regression problem $\min_{x \geq 0} \|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)x - b\|_2$, where $A_i \in \mathbb{R}^{n_i \times d_i}$, $i = 1, \dots, q$ and $b \in \mathbb{R}^{n_1 n_2 \cdots n_q}$. Suppose x is the optimal solution. Similarly, let $S \in \mathbb{R}^{m \times (n_1 n_2 \cdots n_q)}$ be the matrix form of TENSORSKETCH of Section 2. If \tilde{x} is the optimal solution to $\min_{x \geq 0} \|S(A_1 \otimes A_2 \otimes \cdots \otimes A_q)x - Sb\|_2$, we have the following:

Corollary 3.2. (Sketch for tensor nonnegative regression) Suppose $\tilde{x} = \min_{x \geq 0} \|SAx - Sb\|_2$ with TENSORSKETCH $S \in \mathbb{R}^{m \times n}$, where $m = 8(d_1 d_2 \cdots d_q + 1)^2 (2 + 3^q) / (\epsilon^2 \delta)$. Then the following approximation $\|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)\tilde{x} - b\|_2 \leq (1 + \epsilon) \text{OPT}$ holds with probability at least $1 - \delta$, where $\text{OPT} = \min_{x \geq 0} \|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)x - b\|_2$.

The proof of Corollary 3.2 can be found in Appendix C.1.

4 P-SPLINES

B-splines are local basis functions, consisting of low degree (e.g., quadratic, cubic) polynomial segments. The positions where the segments join are called the knots. B-splines have local support and are thus suitable for smoothing and interpolating data with complex patterns. Unfortunately, control over smoothness is limited: one can only change the number and positions of the knots. If there are no reasons to assume that smoothness is non-uniform, the knots will be equally spaced and the only tuning parameter is their (discrete) number. In contrast P-spline (Eilers and Marx,

1996) equally spaces B-splines, discards the derivative completely, and controls smoothness by regularizing the sum of squares of differences of coefficients. Specifically Eilers and Marx Eilers and Marx (1996) proposed the P-spline recipe: (1) use a (quadratic or cubic) B-spline basis with a large number of knots, say 10-50; (2) introduce a penalty on (second or third order) differences of the B-spline coefficients; (3) minimize the resulting penalized likelihood function; (4) tune smoothness with the weight of the penalty, using cross-validation or AIC to determine the optimal weight.

We give a brief overview of B-Splines and P-Splines below. Let b and u , each vectors of length n , represent the observed and explanatory variables, respectively. Once a set of knots is chosen, the B-spline basis A follows from u . If there are d basis functions then A is $n \times d$. In the case of normally distributed observations the model is $b = Ax + e$, with independent errors e . In the case of B-spline regression the sum of squares of residuals $\|b - Ax\|_2$ is minimized and the normal equations $A^T A \hat{x} = A^T b$ are obtained; the explicit solution $\hat{x} = (A^T A)^{-1} A^T b$ results. The P-spline approach minimizes the penalized least-squares function

$$\|b - Ax\|_2^2 + \lambda \|Lx\|_2^2, \quad (1)$$

where $L \in \mathbb{R}^{p \times n}$ is a matrix that forms differences of order ℓ , i.e., $L_\ell x = \Delta^\ell x$. Examples of this matrix, for $\ell = 1$ and $\ell = 2$ are :

$$L_1 = \begin{bmatrix} -1 & 1 & 0 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 0 & -1 & 1 \end{bmatrix}, L_2 = \begin{bmatrix} 1 & -2 & 1 & 0 & 0 \\ 0 & 1 & -2 & 1 & 0 \\ 0 & 0 & 1 & -2 & 1 \end{bmatrix}.$$

The parameter λ determines the influence of the penalty. If λ is zero, we are back to B-spline regression; increasing λ makes \hat{x} , and hence $\hat{b} = A\hat{x}$, smoother.

Let x^* denote $\text{argmin}_{x \in \mathbb{R}^d} \|Ax - b\|_2^2 + \lambda \|Lx\|_2^2$, and OPT denote $\|Ax^* - b\|_2^2 + \lambda \|Lx^*\|_2^2$. In general $x^* = (A^T A + \lambda L^T L)^{-1} A^T b = A^T (AA^T + \lambda LL^T)^{-1} b$, so x^* can be found in $O(\text{nnz}(A) \min(n, d))$ time using an iterative method (e.g., LSQR). Our first goal in this section is to design faster algorithms that find an approximate \tilde{x} in the following sense:

$$\|A\tilde{x} - b\|_2^2 + \lambda \|L\tilde{x}\|_2^2 \leq (1 + \epsilon) \text{OPT}. \quad (2)$$

4.1 Sketching for P-Spline

We first introduce a new definition of *Statistical Dimension* that extends the statistical dimension defined for Ridge Regression (i.e., L is an identity matrix in Eq. 1) (Avron et al., 2016) to P-Spline regression.

The problem (1) can also be analyzed by generalized singular value decomposition (GSVD)

Golub and Van Loan (2013). For matrices $A \in \mathbb{R}^{n \times d}$ and $L \in \mathbb{R}^{p \times d}$ with $\text{rank}(L) = p$ and $\text{rank}\left(\begin{bmatrix} A \\ L \end{bmatrix}\right) = d$, the GSVD of (A, L) is given by the pair of factorizations $A = U \begin{bmatrix} \Sigma & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & I_{d-p} \end{bmatrix} RQ^\top$ and $L = V \begin{bmatrix} \Omega & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & I_{d-p} \end{bmatrix} RQ^\top$, where $U \in \mathbb{R}^{m \times n}$ has orthonormal columns, $V \in \mathbb{R}^{p \times p}$, $Q \in \mathbb{R}^{d \times d}$ are orthogonal, $R \in \mathbb{R}^{d \times d}$ is upper triangular and non-singular, and Σ and Ω are $p \times p$ diagonal matrices: $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_p)$ and $\Omega = \text{diag}(\mu_1, \mu_2, \dots, \mu_p)$ with $0 \leq \sigma_1 \leq \sigma_2 \leq \dots \leq \sigma_p < 1$ and $1 \geq \mu_1 \geq \mu_2 \geq \dots \geq \mu_p > 0$, satisfying $\Sigma^\top \Sigma + \Omega^\top \Omega = I_p$. The *generalized singular values* γ_i of (A, L) are defined by the ratios $\gamma_i = \sigma_i / \mu_i$ ($i = [p]$).

In this section we design an algorithm that is aimed at the case when $n \gg d$. The general strategy is to design a distribution on matrices of size m -by- n (m is a parameter), sample an S from that distribution, and solve $\tilde{x} \equiv \text{argmin}_{x \in \mathbb{R}^d} \|S(Ax - b)\|_2^2 + \lambda \|Lx\|_2^2$.

The following lemma defines conditions on the distribution of S that guarantees Eq. (2) holds with constant probability (which can be boosted to high probability by repetition and taking the minimum objective).

Lemma 4.1. *Let $x^* \in \mathbb{R}^d$, $A \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$ as above. Let $U_1 \in \mathbb{R}^{n \times d}$ denote the first n rows of an orthogonal basis for $\begin{bmatrix} A \\ \sqrt{\lambda}L \end{bmatrix} \in \mathbb{R}^{(n+p) \times d}$. Let sketching matrix $S \in \mathbb{R}^{m \times n}$ have a distribution such that with constant probability*

$$(I) \|U_1^\top S^\top S U_1 - U_1^\top U_1\|_2 \leq 1/4,$$

and

$$(II) \|U_1^\top (S^\top S - I)(b - Ax^*)\|_2 \leq \sqrt{\epsilon \text{OPT}/2}.$$

Let \tilde{x} denote $\text{argmin}_{x \in \mathbb{R}^d} \|S(Ax - b)\|_2^2 + \lambda \|Lx\|_2^2$. Then with probability at least $9/10$,

$$\|A\tilde{x} - b\|_2^2 + \lambda \|L\tilde{x}\|_2^2 \leq (1 + \epsilon) \text{OPT}.$$

Define the *statistical dimension* for P-Splines as follows:

Definition 4.2 (Statistical Dimension for P-Splines). *For S-Spline in Eq. (1), the statistical dimension is defined as $\text{sd}_\lambda(A, L) = \sum_i 1/(1 + \lambda/\gamma_i^2) + d - p$.*

The following theorem shows that there is a sparse subspace embedding matrix $S \in \mathbb{R}^{m \times n}$ (e.g., COUNTSKETCH), with $m \geq K(\text{sd}_\lambda(A, L)/\epsilon + \text{sd}_\lambda(A, L)^2)$, that satisfies Property (I) and (II) of Lemma 4.1, and hence achieves an ϵ -approximation solution to problem 1:

Theorem 4.3. (P-Spline regression) There is a constant $K > 0$ such that for $m \geq K(\epsilon^{-1} \text{sd}_\lambda(A, L) +$

Algorithm 2 P-Spline Tensor product regression

- 1: **procedure** PTRREGRESSION($A, b, K, L, \epsilon, \delta$)
 - 2: $m \leftarrow K(\epsilon^{-1} \text{sd}_\lambda(A, L) + \text{sd}_\lambda(A, L)^2)$
 - 3: Choose S to be a $m \times n$ TENSORSKETCH matrix
 - 4: Compute $S(A_1 \otimes A_2 \otimes \dots \otimes A_q)$ and Sb
 - 5: $\tilde{x} = \text{argmin}_{x \in \mathbb{R}^d} \|S(Ax - b)\|_2^2 + \lambda \|Lx\|_2^2$.
 - 6: **return** \tilde{x}
 - 7: **end procedure**
-

$\text{sd}_\lambda(A, L)^2$) and $S \in \mathbb{R}^{m \times n}$ a sparse embedding matrix (e.g., COUNTSKETCH) with SA computable in $O(\text{nnz}(A))$ time, Property (I) and (II) of Lemma 4.1 apply, and with constant probability the corresponding $\tilde{x} = \text{argmin}_{x \in \mathbb{R}^d} \|S(Ax - b)\|_2^2 + \lambda \|Lx\|_2^2$ is an ϵ -approximate solution to $\min_{x \in \mathbb{R}^d} \|b - Ax\|_2^2 + \lambda \|Lx\|_2^2$.

Note $\text{sd}_\lambda(A, L)$ is upper bounded by d . The above theorem shows that the statistical dimension allows us to design smaller sketch matrices whose size only depends on $O(\text{poly}(\text{sd}_\lambda(A, L)/\epsilon))$ instead of $O(\text{poly}(d/\epsilon))$, without sacrificing the approximation accuracy.

4.2 Tensor Sketching for Multi-Dimensional P-Spline

Tensor products allow a natural extension of one-dimensional P-spline smoothing to multi-dimensional P-Spline. We focus on 2-dimensional P-Spline but our results can be generalized to the multi-dimensional setting. Assume that in addition to u we have a second explanatory variable v . We have data triples $(u_i, v_j, b_{(i-1) \cdot n_2 + j})$ for $i = 1, \dots, n_1$ and $j = 1, \dots, n_2$. We seek a smooth surface $f(u, v)$ which gives a good approximation to the response b . Let $A_1, n \times d_1$, be a B-spline basis along u , and $A_2, n \times d_2$, be a B-spline basis along v . We form the tensor product basis as $A_1 \otimes A_2$. When n_1 and n_2 are large, we do not compute $A_1 \otimes A_2$. We apply TENSORSKETCH here to avoid explicitly forming $A_1 \otimes A_2$ to speed up computation.

Let $X = [x_{kl}]$ be a $d_1 \times d_2$ matrix of coefficients. Then, for given X , the value fit at (u, v) is $f(u, v) = \sum_k \sum_l A_{2,k}(v) A_{1,l}(u) x_{kl}$ and so X may be chosen using least squares by minimizing $\sum_{i,j} [b_{(i-1) \cdot n_2 + j} - f(u_i, v_i)]^2 = \sum_i [b_{(i-1) \cdot n_2 + j} - \sum_k \sum_l A_{2,k}(v_i) A_{1,l}(u_i) x_{kl}]^2$. Using Kronecker product, the above minimization can be written in the form $\min \|b - Ax\|_2$, where $A = A_1 \otimes A_2 \in \mathbb{R}^{n_1 n_2 \times d_1 d_2}$ and $x = \text{vec}(X)$. Again the P-spline approach minimizes the penalized least-squares function $\|b - (A_1 \otimes A_2)x\|_2^2 + \lambda \|Lx\|_2^2$. Consider the tensor p-spline regression problem $\min_x \|(A_1 \otimes A_2 \otimes \dots \otimes A_q)x - b\|_2^2 + \lambda \|Lx\|_2^2$, where $L \in \mathbb{R}^{p \times n}$, $A_i \in \mathbb{R}^{n_i \times d_i}$, $i = 1, \dots, q$ and $b \in \mathbb{R}^n$. Let $S \in \mathbb{R}^{m \times n}$ be the matrix form of TENSORSKETCH of Section 2.

so that SAU is an (α, β, p) -well-conditioned basis of the column space of SA ; 3) Output $\mathcal{AU}/(d\gamma_p)$, where $\gamma_p \equiv \sqrt{2}t^{1/p-1/2}$ for $p \leq 2$, and $\gamma_p \equiv \sqrt{2}w^{1/2-1/p}$ for $p \geq 2$, where $t = \prod_{i=1}^q m_i$ and $w = \prod_{i=1}^q w_i$. The following Lemma 5.3 (the proof can be found in the Appendix) is the analogue of that in Clarkson et al. (2013) proved for the Fast Johnson Lindenstauss Transform.

Lemma 5.3. *For any $p \geq 1$. $\text{Condition}(\mathcal{A})$ computes $\mathcal{AU}/(d\gamma_p)$ which is an $(\alpha, \beta\sqrt{3}d(tw)^{1/p-1/2}, p)$ -well-conditioned basis of \mathcal{A} , with probability at least $1 - \prod_{i=1}^q (n_i/w_i)\delta$.*

Lemma 5.3 indicates that $\text{Condition}(\mathcal{A})$ computes a $(\alpha, \beta \cdot \text{poly}(\max(d, \log n)), p)$ -well-conditioned basis. A well-conditioned basis can be used to solve ℓ_p regression problems, via sampling a subset of rows of the well-conditioned basis \mathcal{AU} with probabilities proportional to the p -th power of the ℓ_p norm of the rows (Woodruff, 2014). However the first issue for sampling is that we cannot afford to compute \mathcal{AU} as this requires $O(\text{nnz}(\mathcal{A})d)$ time. To fix this, we apply a Gaussian sketch matrix $G \in \mathbb{R}^{d \times \log(n)}$ with i.i.d normal random variables (Drineas et al., 2011) on the right hand side of \mathcal{AU} . Note that \mathcal{AUG} can be computed efficiently by first computing UG and then $\mathcal{A}(UG)$ in time $O(d^2 \log(n) + \text{nnz}(\mathcal{A}) \log(n))$. The second issue is that even with \mathcal{AUG} , computing the ℓ_p norm of each row of \mathcal{AUG} takes $O(n)$ time, but we want sublinear time. This leads us to the following sampling technique.

The high level idea is that since \mathcal{AUG} only has $O(\log(n))$ columns, we can afford to sample columns of \mathcal{AUG} with probability proportional to the p -th power of the ℓ_p norms of the columns, if we can efficiently estimate the ℓ_p norms of the columns (note that naively computing the ℓ_p norm of a column also takes $O(n)$ time). Let us denote the first of column of UG as $e \in \mathbb{R}^{\log(n)}$. We focus on how to efficiently estimate the ℓ_p norm of the first column of \mathcal{AUG} , which is $\mathcal{A}e$, and all the left columns can be estimated in the same way. We first reshape the vector $\mathcal{A}e$ into a 2-d matrix

$$M = A_1 \otimes \dots \otimes A_{q_1} E(A_{q_1+1} \otimes \dots \otimes A_q)^\top, \quad (3)$$

where $M \in \mathbb{R}^{(n_1 n_2 \dots n_{q_1}) \times (n_{q_1+1} \dots n_q)}$, E is obtained from reshaping e into a $(d_1 d_2 \dots d_{q_1}) \times (d_{q_1+1} \dots d_q)$ matrix and $q_1 \in [1, q]$ is chosen such that $(n_1 n_2 \dots n_{q_1}) \approx (n_{q_1+1} \dots n_q)$. Namely we reshape the column $\mathcal{A}e$ into a (nearly) square matrix. Focusing now on $p = 1$, note that $\|\mathcal{A}e\|_1 = \sum_{i=1}^{n_{q_1+1} \dots n_q} \|M_{*,i}\|_1$. Hence to estimate $\|\mathcal{A}e\|_1$, we only need to estimate the ℓ_1 norm of the columns of M .

Let us apply a sketch matrix $R \in \mathbb{R}^{O(\log(n)) \times \prod_{i=1}^{q_1} n_i}$, whose entries are sampled i.i.d. from the Cauchy distribution, to the left hand side of M . For the i 'th column of M , let us define random variables $z_i^i = R_{i,*}^\top M_i$,

Algorithm 3 ℓ_1 tensor product regression

- 1: **procedure** L1TRREGRESSION(A, b, ϵ, δ)
 - 2: Construct a TENSORSKETCH $S \in \mathbb{R}^{(m \prod_{i=1}^q \frac{n_i}{w_i}) \times n}$.
 - 3: Run $\text{Condition}(\mathcal{A})$ using S to compute $U/(d\gamma_p)$.
 - 4: Generate a Gaussian matrix $G \in \mathbb{R}^{d \times O(\log(n))}$ and a Cauchy sketch matrix $R \in \mathbb{R}^{\log(n) \times n}$.
 - 5: **for** for each column e in UG **do**
 - 6: Reshape $\mathcal{A}e$ to M (Eq. 3).
 - 7: Compute λ_i and $\lambda_e = \sum_i \lambda_i$ (Eq 4 and 5).
 - 8: **end for**
 - 9: **for** $i \in [\sqrt{\prod_{i=1}^q w_i} \text{poly}(d)]$ **do**
 - 10: Sample a column $(AUG)_{*,e}$ with probability proportional to λ_e . $\triangleright e \in [O(\log(n))]$
 - 11: Reshape $(AUG)_{*,e}$ to M , sample a column $M_{*,j}$ with probability proportional to λ_j .
 - 12: Sample an entry $M_{k,j}$ with probability proportional to $|M_{k,j}|$.
 - 13: Convert (k, j) back to the corresponding row index in \mathcal{A} , denoted as r_i .
 - 14: **end for**
 - 15: $\tilde{x} \leftarrow \min_x \|D(A_1 \otimes A_2 \otimes \dots \otimes A_q)x - Db\|_1$, where D is a diagonal matrix to select the r_1, r_2, \dots, r_N -th rows of \mathcal{A} and b with $N = \sqrt{\prod_{i=1}^q w_i} \text{poly}(d)$.
 - 16: **return** \tilde{x} $\triangleright \tilde{x} \in \mathbb{R}^{d_1 d_2 \dots d_q}$
 - 17: **end procedure**
-

for $l \in [O(\log(n))]$, where $R_{l,*}$ is the l 'th row of R . Due to the 1-stability property of the Cauchy distribution, we have that $\{z_i^i\}_{i=1}^{O(\log(n))}$ are $O(\log(n))$ independent Cauchys scaled by $\|M_i\|_1$. Applying a Chernoff bound to independent half-Cauchys (see Claims 1, 2 and Lemmas 1, 2 in Indyk (2006)), we have $0.5\|M_i\|_1 \leq \text{median}_{l \in [O(\log(n))]} \{|z_l^i|\} \leq 1.5\|M_i\|_1$ with probability at least $1 - 2e^{-cO(\log(n))}$, with constant $c \geq 0.07$. Denote the median of $\{|z_l^i|\}_{l=1}^{O(\log(n))}$ as λ_i . By a union bound over all columns of M , we have that with probability at least $1 - n_{[q] \setminus [q_1]} 2e^{-cO(\log(n))}$:

$$\lambda_i = (1 \pm 0.5)\|M_i\|_1, \forall i \in [n_{[q] \setminus [q_1]}], \quad (4)$$

$$\lambda_e = \sum_{i=1}^{n_{[q] \setminus [q_1]}} \lambda_i = (1 \pm 0.5)\|\mathcal{A}e\|_1 \quad (5)$$

where $n_{[q] \setminus [q_1]} = \prod_{i=q_1+1}^q n_i$.

Note that since \mathcal{AUG} only has $O(\log(n))$ columns, we can afford to compute the ℓ_1 norm of all the columns using the above procedure. Let us denote the ℓ_1 norms of the columns of \mathcal{AUG} by $\lambda_1, \lambda_2, \dots, \lambda_{O(\log(n))}$. We can sample a column $(\mathcal{AUG})_{*,i}$ with probability proportional to λ_i (Line 10 in Alg. 3). Once we sample a column $\mathcal{AUG}_{*,i}$, we need to sample an entry $j \in [n]$

with probability proportional to the absolute value of the entry $|(\mathcal{AUG})_{j,i}|$. As we cannot afford to compute $|(\mathcal{AUG})_{j,i}|$ for all $j \in [n]$, we use the reshaped 2-d matrix M of $(\mathcal{AUG})_{*,i}$. Note that sampling an entry in M with probability proportional to the absolute values of entries of M is equivalent to sampling an entry $j \in [n]$ from $(\mathcal{AUG})_{*,i}$ with probability proportional to the absolute value of the entries of $(\mathcal{AUG})_{*,i}$. We first sample a column $M_{*,j}$ from all the columns of M with probability proportional to λ_j for $j \in [\prod_{i=q_1+1}^q n_i]$. We then sample an entry $M_{k,j}$ with probability proportional to $|M_{k,j}|$ for $k \in [\prod_{i=1}^{q_1} n_i]$. Noting that $k \in [\prod_{i=1}^{q_1} n_i]$ and $j \in [\prod_{i=q_1+1}^q n_i]$, the pair (k, j) uniquely determines a corresponding row index r in \mathcal{A} , for some $r \in [\prod_{i=1}^q n_i]$. Hence we successfully sample a row from \mathcal{AUG} without ever computing the ℓ_p norm of the rows. The above sampling procedure is summarized in Line 10 to Line 13 in Alg. 3. We use the above procedure to sample $\sqrt{\prod_{i=1}^q w_i} \text{poly}(d)$ rows of \mathcal{A} . Let D be a diagonal matrix that selects the corresponding sampled rows from \mathcal{A} . We can now solve a smaller ADL problem as $\min_x \|D\mathcal{A} - Db\|_1$. Note that our analysis focuses on the ℓ_1 norm. We can extend the analysis to general ℓ_p norms by using a sketching matrix $R \in \mathbb{R}^{O(\log n) \times \prod_{i=1}^{q_1} n_i}$ with entries sampled i.i.d. from a p -stable distribution for $p \in [1, 2]$.

We now present our main theorem and defer the proofs to the appendix:

Theorem 5.4. (Main result) Given $\epsilon \in (0, 1)$, $\mathcal{A} \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$, Alg. 3 computes \hat{x} such that with probability at least $1/2$, $\|\mathcal{A}\hat{x} - b\|_1 \leq (1 + \epsilon) \min_{x \in \mathbb{R}^d} \|\mathcal{A}x - b\|_1$. For the special case when $q = 2$, $n_1 = n_2$, the algorithm's running time is $O(n_1^{3/2} \text{poly}(\prod_{i=1}^2 d_i/\epsilon))$.

For the special case where $q = 2$ and $n_1 = n_2$, we can see from theorem 5.4 our algorithm computes \hat{x} in time $O(n_1^{3/2} \text{poly}(\hat{d}))$, which is faster than $O(n_1^2)$ —the time needed for forming $A_1 \otimes A_2$. Note that we can run Alg. 3 $O(\log(1/\delta))$ times independently and pick the best solution among these independent runs to boost the success probability to be $1 - \delta$, for $\delta \in [0, 1)$.

6 NUMERICAL EXPERIMENTS

We generate matrices A_1 , A_2 and b with all entries sampled i.i.d from a normal distribution. The baseline we compared to is directly solving regression without sketching. We let T_1 be the time for directly solving the regression problem, and T_2 be the time of our algorithm. The time ratio is $r_t = T_2/T_1$. The relative residual percentage is defined by $r_e = \frac{100 \|(A_1 \otimes A_2)\tilde{x} - b\|_2 - \|(A_1 \otimes A_2)x^* - b\|_2}{\|(A_1 \otimes A_2)x^* - b\|_2}$, where \tilde{x} is the output of our algorithms and x^* is the optimal solution.

Table 1: Examples 6.1 and 6.2: the values of r_e and r_t with respect to different sampling parameters m .

	m	r_e	r_t		m	r_e	r_t
ℓ_2	8000	1.79%	0.11	ℓ_1	8000	1.89%	0.06
	12000	1.24%	0.18		12000	1.33%	0.11
	16000	1.01%	0.25		16000	0.992%	0.18

Throughout the simulations, we use a moderate input matrix size in order to accommodate the brute force algorithm and to compare to the exact solution.

Example 6.1 (ℓ_2 Regression). We create a design matrix with moderate size by fixing $n_1 = n_2 = 300$ and $d_1 = d_2 = 15$. Thus $\mathcal{A} \in \mathbb{R}^{90000 \times 225}$. We do 10 rounds to compute the mean values of r_e and r_t , which is reported in Table 1 (Left).

From Table 1 (left), we can see that when the number m of sampled rows in Algorithm 1 increases from 8000 to 12000, the mean values of r_e decrease while the mean values of r_t increase. In general we can see that we can achieve around 1% relative error while being 5 times faster than the direct method.

Example 6.2 (ℓ_1 Regression). We set $n_1 = n_2 = 300$, $d_1 = d_2 = 15$. For $\min_x \|Ax - b\|_1$, we solve it by a Linear Programming solver in Gurobi (Gurobi Optimization, 2016). We tested different numbers of sampled rows m (the number of rows in D in Alg. 3). The results are summarized in Table 1 (Right).

As we can see from Table 1 (right), our method is around 10 times faster than directly solving the problem, with relative error only around 1%.

Example 6.3 (P-Spline Regression). For P-spline regression, L_3 is fixed. We use 30 knots and cubic B-splines. The data $u = (u_i) \in \mathbb{R}^{n_1}$, $v = (v_j) \in \mathbb{R}^{n_2}$ and $b \in \mathbb{R}^{n_1 n_2}$ are generated i.i.d from the normal distribution. We compute the B-spline basis matrices $A_1 \in \mathbb{R}^{n_1 \times d_1}$ and $A_2 \in \mathbb{R}^{n_2 \times d_2}$ separately, if there are d_1 and d_2 basis functions for the B-spline. We set $n_1 = n_2$, $d_1 = d_2$, with $n_1 n_2 = 10^4$ and $d_1 d_2 = 529$. The original and sketched P-spline regression problem are both solved by REGULARIZATION TOOLS (Hansen, 1994) via computing their GSVDs. We test different choices of λ . The results are shown in Table 2.

From Table 2, sampling only 20% of the rows can give around 0.05% relative error, while the computation time is half of the time of directly solving the p-spline.

7 CONCLUSION

We propose algorithms for efficiently solving tensor product least squares regression, regularized P-splines, as well as tensor product least absolute deviation regression, using sketching techniques. Our main contributions are: (1) we apply TENSORSKETCH to least

Table 2: Example 6.3: the mean values of r_e and r_t with respect to different sampling parameters m and regularization parameters λ .

λ	m	r_e	r_t
1	2000	4.43e-2%	0.52
	4000	2.99e-2%	0.70
	6000	7.92e-2%	0.91
0.1	2000	6.98e-2%	0.47
	4000	4.07e-2%	0.72
	6000	5.41e-2%	0.94
0.01	2000	3.78e-2%	0.46
	4000	1.07e-1%	0.71
	6000	2.97e-2%	0.95

square regression problems, (2) we propose new statistical dimension measures for P-splines, extending the previous statistical dimension defined only for classic Ridge regression, and (3) we extend TENSORSKETCH to ℓ_p norms and propose an algorithm that can solve tensor product ℓ_1 regression in time sublinear in the time for explicitly computing the tensor product. Simulation results support our theorems and demonstrate that our algorithms are much faster than brute-force algorithms and can achieve approximate solutions that are close to optimal.

ACKNOWLEDGEMENT

Huaian Diao is supported in part by the Fundamental Research Funds for the Central Universities under the grant 2412017FZ007. Wen Sun is supported in part by Office of Naval Research contract N000141512365.

References

- H. Avron, H. Nguyen, and D. Woodruff. Subspace embeddings for the polynomial kernel. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2258–2266, 2014.
- H. Avron, K. L. Clarkson, and D. P. Woodruff. Sharper bounds for regression and low-rank approximation with regularization. *CoRR*, abs/1611.03225, 2016.
- L. Carter and M. N. Wegman. Universal classes of hash functions. *J. Comput. Syst. Sci.*, 18(2):143–154, 1979.
- M. Charikar, K. Chen, and M. Farach-Colton. Finding frequent items in data streams. *Theor. Comput. Sci.*, 312(1):3–15, 2004.
- D. Chen and R. J. Plemmons. Nonnegativity constraints in numerical analysis. In *The birth of numerical analysis*, pages 109–139. World Sci. Publ., Hackensack, NJ, 2010.
- K. Clarkson, P. Drineas, M. Magdon-Ismail, M. Mahoney, X. Meng, and D. P. Woodruff. The fast Cauchy transform and faster robust linear regression. In *SODA*, 2013.
- K. L. Clarkson. Subgradient and sampling algorithms for ℓ_1 regression. In *Proceedings of the sixteenth annual ACM-SIAM symposium on Discrete algorithms (SODA)*, pages 257–266, 2005.
- K. L. Clarkson and D. 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. <https://arxiv.org/pdf/1207.6365>, 2013.
- A. Dasgupta, P. Drineas, B. Harb, R. Kumar, and M. W. Mahoney. Sampling algorithms and coresets for ℓ_p regression. *SIAM J. Comput.*, 38(5):2060–2078, 2009.
- P. Drineas, M. Magdon-Ismail, M. W. Mahoney, and D. P. Woodruff. Fast approximation of matrix coherence and statistical leverage. *CoRR*, abs/1109.3843, 2011.
- P. H. Eilers and B. D. Marx. Multidimensional density smoothing with p-splines. In *Proceedings of the 21st international workshop on statistical modelling*, 2006.
- P. H. C. Eilers and B. D. Marx. Flexible smoothing with B-splines and penalties. *Statist. Sci.*, 11(2): 89–121, 1996.
- P. H. C. Eilers, B. D. Marx, and M. Durbán. Twenty years of P-splines. *SORT*, 39(2):149–186, 2015. ISSN 1696-2281.
- D. W. Fausett and C. T. Fulton. Large least squares problems involving Kronecker products. *SIAM J. Matrix Anal. Appl.*, 15(1):219–227, 1994.
- G. H. Golub and C. F. Van Loan. *Matrix computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, MD, 2013.
- I. Gurobi Optimization. Gurobi optimizer reference manual, 2016. URL <http://www.gurobi.com>.
- P. C. Hansen. Regularization tools: A matlab package for analysis and solution of discrete ill-posed problems. *Numerical Algorithms*, 6(1):1–35, 1994.
- P. Indyk. Stable distributions, pseudorandom generators, embeddings, and data stream computation. *J. ACM*, 53(3):307–323, 2006.
- R. Kannan and S. Vempala. Spectral algorithms. *Foundations and Trends in Theoretical Computer Science*, 4(3-4):157–288, 2009.

- Y. Liang, M.-F. F. Balcan, V. Kanchanapally, and D. Woodruff. Improved distributed principal component analysis. In *Advances in Neural Information Processing Systems*, pages 3113–3121, 2014.
- M. W. Mahoney. Randomized algorithms for matrices and data. *Foundations and Trends in Machine Learning*, 3(2):123–224, 2011.
- X. Meng and M. W. Mahoney. Low-distortion subspace embeddings in input-sparsity time and applications to robust linear regression. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 91–100. ACM, <https://arxiv.org/pdf/1210.3135>, 2013.
- J. Nelson and H. L. Nguyen. Osnap: Faster numerical linear algebra algorithms via sparser subspace embeddings. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 117–126. IEEE, <https://arxiv.org/pdf/1211.1002>, 2013.
- S. Oh, S. Kwon and J. Yun. A method for structured linear total least norm on blind deconvolution problem. *Applied Mathematics and Computing*, 19:151–164, 2005.
- R. Pagh. Compressed matrix multiplication. *ACM Trans. Comput. Theory*, 5(3):9:1–9:17, 2013.
- M. Patrascu and M. Thorup. The power of simple tabulation hashing. *J. ACM*, 59(3):14, 2012.
- N. Pham and R. Pagh. Fast and scalable polynomial kernels via explicit feature maps. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining(KDD)*, pages 239–247. ACM, 2013.
- P. J. Rousseeuw and A. M. Leroy. *Robust regression and outlier detection*, volume 589. John wiley & sons, 2005.
- C. Sohler and D. P. Woodruff. Subspace embeddings for the ℓ_1 -norm with applications. In *Proceedings of the forty-third annual ACM symposium on Theory of computing (STOC)*, pages 755–764. ACM, 2011.
- Z. Song, D. P. Woodruff, and P. Zhong. Low rank approximation with entrywise ℓ_1 -norm error. In *Proceedings of the 49th Annual Symposium on the Theory of Computing (STOC)*. ACM, <https://arxiv.org/pdf/1611.00898>, 2017a.
- Z. Song, D. P. Woodruff, and P. Zhong. Relative error tensor low rank approximation. *arXiv preprint arXiv:1704.08246*, 2017b.
- C. Van Loan. *Computational frameworks for the fast Fourier transform*, volume 10 of *Frontiers in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1992.
- C. F. Van Loan and N. Pitsianis. Approximation with Kronecker products. In *Linear algebra for large scale and real-time applications (Lewven, 1992)*, volume 232 of *NATO Adv. Sci. Inst. Ser. E Appl. Sci.*, pages 293–314. Kluwer Acad. Publ., Dordrecht, 1993.
- D. P. Woodruff. Sketching as a tool for numerical linear algebra. *Foundations and Trends in Theoretical Computer Science*, 10(1-2):1–157, 2014.
- Y. Zheng, G. Liu, S. Sugimoto, S. Yan, and M. Okutomi. Practical low-rank matrix approximation under robust ℓ_1 -norm. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 1410–1417. IEEE, 2012.

A Background: COUNTSKETCH and TENSORSKETCH

We start by describing the COUNTSKETCH transform [Charikar et al. \(2004\)](#). Let m be the target dimension. When applied to n -dimensional vectors, the transform is specified by a 2-wise independent hash function $h : [n] \rightarrow [m]$ and a 2-wise independent sign function $s : [n] \rightarrow \{-1, +1\}$. When applied to v , the value at coordinate i of the output, $i = 1, 2, \dots, m$ is $\sum_{j|h(j)=i} s(j)v_j$. Note that COUNTSKETCH can be represented as an $m \times n$ matrix in which the j -th column contains a single non-zero entry $s(j)$ in the $h(j)$ -th row.

We now describe the TENSORSKETCH transform [Pagh \(2013\)](#). Suppose we are given points $v_i \in \mathbb{R}^{n_i}$, where $i = 1, \dots, q$ and so $\phi(v_1, \dots, v_q) = v_1 \otimes v_2 \otimes \dots \otimes v_q \in \mathbb{R}^{n_1 n_2 \dots n_q}$, and the target dimension is again m . The transform is specified using q 3-wise independent hash functions $h_i : [n_i] \rightarrow [m]$, and q 4-wise independent sign functions $s_i : [n_i] \rightarrow \{+1, -1\}$, where $i = 1, \dots, q$. TENSORSKETCH applied to v_1, \dots, v_q is then COUNTSKETCH applied to $\phi(v_1, \dots, v_q)$ with hash function $H : [n_1 n_2 \dots n_q] \rightarrow [m]$ and sign function $S : [n_1 n_2 \dots n_q] \rightarrow \{+1, -1\}$ defined as follows:

$$H(i_1, \dots, i_q) = h_1(i_1) + h_2(i_2) + \dots + h_q(i_q) \bmod m,$$

and

$$S(i_1, \dots, i_q) = s_1(i_1) \cdot s_2(i_2) \cdot \dots \cdot s_q(i_q),$$

where $i_j \in [n_j]$. It is well-known that if H is constructed this way, then it is 3-wise independent [Carter and Wegman \(1979\)](#); [Patrascu and Thorup \(2012\)](#). Unlike the work of Pham and Pagh [Pham and Pagh \(2013\)](#), which only used that H was 2-wise independent, our analysis needs this stronger property of H .

The TENSORSKETCH transform can be applied to v_1, \dots, v_q without computing $\phi(v_1, \dots, v_q)$ as follows. Let $v_j = (v_{j\ell}) \in \mathbb{R}^{n_j}$. First, compute the polynomials

$$p_\ell(x) = \sum_{i=0}^{B-1} x^i \sum_{j_\ell|h_\ell(j_\ell)=i} v_{j_\ell} \cdot s_\ell(j_\ell),$$

for $\ell = 1, 2, \dots, q$. A calculation [Pagh \(2013\)](#) shows

$$\prod_{\ell=1}^q p_\ell(x) \bmod (x^B - 1) = \sum_{i=0}^{B-1} x^i \sum_{(j_1, \dots, j_q)|H(j_1, \dots, j_q)=i} v_{j_1} \cdot \dots \cdot v_{j_q} S(j_1, \dots, j_q),$$

that is, the coefficients of the product of the q polynomials mod $(x^B - 1)$ form the value of TENSORSKETCH(v_1, \dots, v_q). Pagh observed that this product of polynomials can be computed in $O(qm \log m)$ time using the Fast Fourier Transform. As it takes $O(q \max(\text{nnz}(v_i)))$ time to form the q polynomials, the overall time to compute TENSORSKETCH(v) is $O(q(\max(\text{nnz}(v_i)) + m \log m))$.

B TENSORSKETCH is an Oblivious Subspace Embedding (OSE)

Let S be the $m \times (n_1 n_2 \dots n_q)$ matrix such that TENSORSKETCH(v_1, \dots, v_q) is $S \cdot \phi(v_1, \dots, v_q)$ for a randomly selected TENSORSKETCH. Notice that S is a random matrix. In the rest of the paper, we refer to such a matrix as a TENSORSKETCH matrix with an appropriate number of rows, i.e., the number of hash buckets. We will show that S is an oblivious subspace embedding for subspaces in $\mathbb{R}^{n_1 n_2 \dots n_q}$ for appropriate values of m . Notice that S has exactly one non-zero entry per column. The index of the non-zero in the column (i_1, \dots, i_q) is $H(i_1, \dots, i_q) = \sum_{j=1}^q h_j(i_j) \bmod m$. Let $\delta_{a,b}$ be the indicator random variable of whether $S_{a,b}$ is non-zero. The sign of the non-zero entry in column (i_1, \dots, i_q) is $S(i_1, \dots, i_q) = \prod_{j=1}^q s_j(i_j)$. We show that the embedding matrix S of TENSORSKETCH can be used to approximate matrix product and is an oblivious subspace embedding (OSE).

Theorem B.1. *Let S be the $m \times (n_1 n_2 \dots n_q)$ matrix such that*

$$\text{TENSORSKETCH}(v_1, \dots, v_q)$$

is $S \cdot \phi(v_1, \dots, v_q)$ for a randomly selected TENSORSKETCH. The matrix S satisfies the following two properties.

1. (Approximate Matrix Product :) Let A and B be matrices with $n_1 n_2 \cdots n_q$ rows. For $m \geq (2 + 3^q)/(\epsilon^2 \delta)$, we have

$$\Pr_S [\|A^\top S^\top SB - A^\top B\|_F^2 \leq \epsilon^2 \|A\|_F^2 \|B\|_F^2] \geq 1 - \delta.$$

2. (Subspace Embedding :) Consider a fixed k -dimensional subspace V . If $m \geq k^2(2 + 3^q)/(\epsilon^2 \delta)$, then with probability at least $1 - \delta$, $\|Sx\| = (1 \pm \epsilon)\|x\|$ simultaneously for all $x \in V$.

We establish the theorem via two lemmas as in Avron et al. (2016). The first lemma proves the approximate matrix product property via a careful second moment analysis.

Lemma B.2 (Approximate matrix product). Let A and B be matrices with $n_1 n_2 \cdots n_q$ rows. For $m \geq (2 + 3^q)/(\epsilon^2 \delta)$, we have

$$\Pr_S [\|A^\top S^\top SB - A^\top B\|_F^2 \leq \epsilon^2 \|A\|_F^2 \|B\|_F^2] \geq 1 - \delta.$$

Proof. The proof follows that in Avron et al. (2016). Let $C = A^\top S^\top SB$. We have

$$C_{u,u'} = \sum_{t=1}^m \sum_{i,j \in [n_1 n_2 \cdots n_q]} S(i)S(j) \delta_{t,i} \delta_{t,j} A_{i,u} B_{j,u'} = \sum_{t=1}^m \sum_{i \neq j \in [n_1 n_2 \cdots n_q]} S(i)S(j) \delta_{t,i} \delta_{t,i} A_{i,u} B_{j,u'} + (A^\top B)_{u,u'}$$

Thus, $\mathbf{E}[C_{u,u'}] = (A^\top B)_{u,u'}$.

Next, we analyze $\mathbf{E}[(C - A^\top B)_{u,u'}^2]$. We have

$$((C - A^\top B)_{u,u'})^2 = \sum_{t_1, t_2=1}^m \sum_{i_1 \neq j_1, i_2 \neq j_2 \in [n_1 n_2 \cdots n_q]} S(i_1)S(i_2)S(j_1)S(j_2) \cdot \delta_{t_1, i_1} \delta_{t_1, j_1} \delta_{t_2, i_2} \delta_{t_2, j_2} \cdot A_{i_1, u} A_{i_2, u} B_{j_1, u'} B_{j_2, u'}$$

For a term in the summation on the right hand side to have a non-zero expectation, it must be the case that $\mathbf{E}[S(i_1)S(i_2)S(j_1)S(j_2)] \neq 0$. Note that $S(i_1)S(i_2)S(j_1)S(j_2)$ is a product of random signs (possibly with multiplicities) where the random signs in different coordinates in $\{1, \dots, q\}$ are independent and they are 4-wise independent within each coordinate. Thus, $\mathbf{E}[S(i_1)S(i_2)S(j_1)S(j_2)]$ is either 1 or 0. For the expectation to be 1, all random signs must appear with even multiplicities. In other words, in each of the q coordinates, the 4 coordinates of i_1, i_2, j_1, j_2 must be the same number appearing 4 times or 2 distinct numbers, each appearing twice. All the subsequent claims in the proof regarding i_1, i_2, j_1, j_2 agreeing on some coordinates follow from this property.

Let S_1 be the set of coordinates where i_1 and i_2 agree. Note that j_1 and j_2 must also agree in all coordinates in S_1 by the above argument. Let $S_2 \subset [q] \setminus S_1$ be the coordinates among the remaining where i_1 and j_1 agree. Finally, let $S_3 = [q] \setminus (S_1 \cup S_2)$. All coordinates in S_3 of i_1 and j_2 must agree. Similarly as before, note that i_2 and j_2 agree on all coordinates in S_2 and i_2 and j_1 agree on all coordinates in S_3 . We can rewrite $i_1 = (a, b, c), i_2 = (a, e, f), j_1 = (g, b, f), j_2 = (g, e, c)$ where $a = (a_\ell), g = (g_\ell)$ with $\ell \in S_1, b = (b_\ell), e = (e_\ell)$ with $\ell \in S_2$ and $c = (c_\ell), f = (f_\ell)$ with $\ell \in S_3$.

First we show that the contribution of the terms where $i_1 = i_2$ or $i_1 = j_2$ is bounded by $\frac{2\|A_u\|_2^2 \|B_{u'}\|_2^2}{m}$, where A_u is the u th column of A and $B_{u'}$ is the u' th column of B . Indeed, consider the case $i_1 = i_2$. As observed before, we must have $j_1 = j_2$ to get a non-zero contribution. Note that if $t_1 \neq t_2$, we always have $\delta_{t_1, i_1} \delta_{t_2, i_2} = 0$ as $H(i_1)$ cannot be equal to both t_1 and t_2 . Thus, for fixed $i_1 = i_2, j_1 = j_2$,

$$\begin{aligned} & \mathbf{E} \left[\sum_{t_1, t_2=1}^m S(i_1)S(i_2)S(j_1)S(j_2) \cdot \delta_{t_1, i_1} \delta_{t_1, j_1} \delta_{t_2, i_2} \delta_{t_2, j_2} \cdot A_{i_1, u} A_{i_2, u} B_{j_1, u'} B_{j_2, u'} \right] \\ &= \mathbf{E} \left[\sum_{t_1=1}^m \delta_{i_1, t_1}^2 \delta_{j_1, t_1}^2 A_{i_1, u}^2 B_{j_1, u'}^2 \right] \\ &= \frac{A_{i_1, u}^2 B_{j_1, u'}^2}{m} \end{aligned}$$

Summing over all possible values of i_1, j_1 , we get the desired bound of $\frac{\|A_u\|_2^2 \|B_{u'}\|_2^2}{m}$. The case $i_1 = j_2$ is analogous.

Next we compute the contribution of the terms where $i_1 \neq i_2, j_1, j_2$ i.e., there are at least 3 distinct numbers among i_1, i_2, j_1, j_2 . Notice that $\mathbf{E}[\delta_{t_1, i_1} \delta_{t_1, j_1} \delta_{t_2, i_2} \delta_{t_2, j_2}] \leq \frac{1}{m^3}$ because the $\delta_{t, i}$'s are 3-wise independent. For fixed i_1, j_1, i_2, j_2 , there are m^2 choices of t_1, t_2 so the total contribution to the expectation from terms with the same i_1, j_1, i_2, j_2 is bounded by $m^2 \cdot \frac{1}{m^3} \cdot |A_{i_1, u} A_{i_2, u} B_{j_1, u'} B_{j_2, u'}| = \frac{1}{m} |A_{i_1, u} A_{i_2, u} B_{j_1, u'} B_{j_2, u'}|$.

Therefore,

$$\begin{aligned}
 & \mathbf{E}[\|(C - A^\top B)_{u, u'}\|^2] \\
 & \leq \frac{2\|A_u\|_2^2 \|B_{u'}\|_2^2}{m} + \frac{1}{m} \sum_{\text{partition } S_1, S_2, S_3} \sum_{a, g, b, e, c, f} |A_{(a, b, c), u} B_{(g, b, f), u'} A_{(a, e, f), u} B_{(g, e, c), u'}| \\
 & \leq \frac{2\|A_u\|_2^2 \|B_{u'}\|_2^2}{m} + \frac{3^q}{m} \sum_{a, b, c, g, e, f} |A_{(a, b, c), u} B_{(g, b, f), u'} A_{(a, e, f), u} B_{(g, e, c), u'}| \\
 & \leq \frac{2\|A_u\|_2^2 \|B_{u'}\|_2^2}{m} + \frac{3^q}{m} \sum_{g, e, f} \left(\sum_{a, b, c} A_{(a, b, c), u}^2 \right)^{1/2} \left(\sum_{a, b, c} B_{(g, b, f), u'}^2 A_{(a, e, f), u}^2 B_{(g, e, c), u'}^2 \right)^{1/2} \\
 & = \frac{2\|A_u\|_2^2 \|B_{u'}\|_2^2}{m} + \frac{3^q \|A_u\|}{m} \sum_{g, e, f} \left(\sum_b B_{(g, b, f), u'}^2 \right)^{1/2} \left(\sum_{a, c} A_{(a, e, f), u}^2 B_{(g, e, c), u'}^2 \right)^{1/2} \\
 & \leq \frac{2\|A_u\|_2^2 \|B_{u'}\|_2^2}{m} + \frac{3^q \|A_u\|}{m} \sum_e \left(\sum_{b, g, f} B_{(g, b, f), u'}^2 \right)^{1/2} \left(\sum_{a, c, g, f} A_{(a, e, f), u}^2 B_{(g, e, c), u'}^2 \right)^{1/2} \\
 & = \frac{2\|A_u\|_2^2 \|B_{u'}\|_2^2}{m} + \frac{3^q \|A_u\| \cdot \|B_{u'}\|}{m} \sum_e \left(\sum_{a, f} A_{(a, e, f), u}^2 \right)^{1/2} \left(\sum_{g, c} B_{(g, e, c), u'}^2 \right)^{1/2} \\
 & \leq \frac{2\|A_u\|_2^2 \|B_{u'}\|_2^2}{m} + \frac{3^q \|A_u\| \cdot \|B_{u'}\|}{m} \left(\sum_{a, e, f} A_{(a, e, f), u}^2 \right)^{1/2} \left(\sum_{g, e, c} B_{(g, e, c), u'}^2 \right)^{1/2} \\
 & = \frac{(2 + 3^q) \|A_u\|_2^2 \|B_{u'}\|_2^2}{m},
 \end{aligned}$$

where the second inequality follows from the fact that there are at most 3^q partitions of $[q]$ into 3 sets. The other inequalities are from Cauchy-Schwarz.

Combining the above bounds, we have $\mathbf{E}[\|(C - A^\top B)_{u, u'}\|^2] \leq \frac{(2+3^q)\|A_u\|_2^2 \|B_{u'}\|_2^2}{m}$. For $m \geq (2 + 3^q)/(\epsilon^2 \delta)$, by the Markov inequality, $\|A^\top S^\top SB - A^\top B\|_F^2 \leq \epsilon^2 \|A\|_F^2 \|B\|_F^2$ with probability $1 - \delta$. \square \square

The second lemma proves that the subspace embedding property follows from the approximate matrix product property.

Lemma B.3 (Oblivious subspace embeddings). *Consider a fixed k -dimensional subspace $V \subset \mathbb{R}^{n_1 n_2 \cdots n_q}$. If $m \geq k^2(2 + 3^q)/(\epsilon^2 \delta)$, then with probability at least $1 - \delta$, $\|Sx\|_2 = (1 \pm \epsilon)\|x\|_2$ simultaneously for all $x \in V$.*

Proof. Let B be a $(n_1 n_2 \cdots n_q) \times k$ matrix whose columns form an orthonormal basis of V . Thus, we have $B^\top B = I_k$ and $\|B\|_F^2 = k$. The condition that $\|Sx\|_2 = (1 \pm \epsilon)\|x\|_2$ simultaneously for all $x \in V$ is equivalent to the condition that the singular values of SB are bounded by $1 \pm \epsilon$. By Lemma B.2, for $m \geq (2 + 3^q)/((\epsilon/k)^2 \delta)$, with probability at least $1 - \delta$, we have

$$\|B^\top S^\top SB - B^\top B\|_F^2 \leq (\epsilon/k)^2 \|B\|_F^4 = \epsilon^2$$

Thus, we have $\|B^\top S^\top SB - I_k\|_2 \leq \|B^\top S^\top SB - I_k\|_F \leq \epsilon$. In other words, the squared singular values of SB are bounded by $1 \pm \epsilon$, implying that the singular values of SB are also bounded by $1 \pm \epsilon$. Note that $\|A\|_2$ for a matrix A denotes its operator norm. \square \square

C Missing Proofs

C.1 Proofs for Tensor Product Least Square Regression

Theorem 3.1. (Tensor regression) Suppose \tilde{x} is the output of Algorithm 1 with TENSORSKETCH $S \in \mathbb{R}^{m \times n}$, where $m = 8(d_1 d_2 \cdots d_q + 1)^2(2 + 3^q)/(\epsilon^2 \delta)$. Then the following approximation $\|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)\tilde{x} - b\|_2 \leq (1 + \epsilon) \text{OPT}$, holds with probability at least $1 - \delta$.

Proof. It is easy to see that

$$\|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)x - b\|_2 = \left\| \begin{bmatrix} (A_1 \otimes A_2 \otimes \cdots \otimes A_q) & b \end{bmatrix} \begin{bmatrix} x \\ -1 \end{bmatrix} \right\|_2,$$

and identifying

$$y = \begin{bmatrix} (A_1 \otimes A_2 \otimes \cdots \otimes A_q) & b \end{bmatrix} \begin{bmatrix} x \\ -1 \end{bmatrix} \in \mathbb{R}^{n_1 n_2 \cdots n_q}$$

and y is a vector of a subspace $V \subset \mathbb{R}^{n_1 n_2 \cdots n_q}$ with dimension at most $d_1 d_2 \cdots d_q + 1$, we can use Lemma B.3 to conclude that

$$\Pr[\|Sy\|_2 - \|y\|_2 \leq \epsilon\|y\|_2] \geq 1 - \delta$$

when $m = (d_1 d_2 \cdots d_q + 1)^2(2 + 3^q)/(\epsilon^2 \delta)$.

Thus we have

$$\|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)\tilde{x} - b\|_2 \leq \frac{1}{1 - \epsilon} \|S(A_1 \otimes A_2 \otimes \cdots \otimes A_q)\tilde{x} - Sb\|_2$$

and

$$\|S(A_1 \otimes A_2 \otimes \cdots \otimes A_q)x - Sb\|_2 \leq (1 + \epsilon)\|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)x - b\|_2$$

hold with probability at least $1 - \delta$. Then using a union bound, we have

$$\begin{aligned} & \|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)\tilde{x} - b\|_2 \\ & \leq \frac{1}{1 - \epsilon} \|S(A_1 \otimes A_2 \otimes \cdots \otimes A_q)\tilde{x} - Sb\|_2 \\ & \leq \frac{1}{1 - \epsilon} \|S(A_1 \otimes A_2 \otimes \cdots \otimes A_q)x - Sb\|_2 \\ & \leq \frac{1 + \epsilon}{1 - \epsilon} \|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)x - b\|_2 \end{aligned}$$

holds with probability at least $1 - 2\delta$. □

Corollary 3.2. (Sketch for tensor nonnegative regression) Suppose $\tilde{x} = \min_{x \geq 0} \|SAx - Sb\|_2$ with TENSORSKETCH $S \in \mathbb{R}^{m \times n}$, where $m = 8(d_1 d_2 \cdots d_q + 1)^2(2 + 3^q)/(\epsilon^2 \delta)$. Then the following approximation $\|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)\tilde{x} - b\|_2 \leq (1 + \epsilon) \text{OPT}$ holds with probability at least $1 - \delta$, where $\text{OPT} = \min_{x \geq 0} \|(A_1 \otimes A_2 \otimes \cdots \otimes A_q)x - b\|_2$.

Proof. The proof of Theorem 3.2 is similar to the proof of theorem 3.1. Denote $\tilde{x} = \min_{x \geq 0} \|SAx - Sb\|_2$ and $x^* = \min_{x \geq 0} \|Ax - b\|_2$. Using Lemma B.3, we have:

$$\|A\tilde{x} - b\|_2 \leq \frac{1}{1 - \epsilon} \|SA\tilde{x} - Sb\|_2, \tag{6}$$

with probability at least $1 - \delta$, and

$$\|SAx^* - Sb\|_2 \leq (1 + \epsilon)\|Ax^* - b\|_2, \tag{7}$$

with probability at least $1 - \delta$. Hence applying a union bound we have:

$$\begin{aligned} & \|A\tilde{x} - b\|_2 \\ & \leq \frac{1}{1-\epsilon} \|SA\tilde{x} - Sb\|_2 \end{aligned} \tag{8}$$

$$\begin{aligned} & \leq \frac{1}{1-\epsilon} \|SAx^* - Sb\|_2 \\ & \leq \frac{1+\epsilon}{1-\epsilon} \|Ax^* - b\|_2, \end{aligned} \tag{9}$$

with probability at least $1 - 2\delta$. \square

C.2 Proofs for P-Splines

Lemma 4.1. *Let $x^* \in \mathbb{R}^d$, $A \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$ as above. Let $U_1 \in \mathbb{R}^{n \times d}$ denote the first n rows of an orthogonal basis for $\begin{bmatrix} A \\ \sqrt{\lambda}L \end{bmatrix} \in \mathbb{R}^{(n+p) \times d}$. Let sketching matrix $S \in \mathbb{R}^{m \times n}$ have a distribution such that with constant probability*

$$(I) \|U_1^\top S^\top S U_1 - U_1^\top U_1\|_2 \leq 1/4,$$

and

$$(II) \|U_1^\top (S^\top S - I)(b - Ax^*)\|_2 \leq \sqrt{\epsilon \text{OPT}/2}.$$

Let \tilde{x} denote $\operatorname{argmin}_{x \in \mathbb{R}^d} \|S(Ax - b)\|_2^2 + \lambda \|Lx\|_2^2$. Then with probability at least $9/10$,

$$\|A\tilde{x} - b\|_2^2 + \lambda \|L\tilde{x}\|_2^2 \leq (1 + \epsilon) \text{OPT}.$$

Proof. Let $\hat{A} \in \mathbb{R}^{(n+d) \times d}$ have orthonormal columns with $\operatorname{range}(\hat{A}) = \operatorname{range}\left(\begin{bmatrix} A \\ \sqrt{\lambda}L \end{bmatrix}\right)$. (An explicit expression for one such \hat{A} is given below.) Let $\hat{b} \equiv \begin{bmatrix} b \\ 0_d \end{bmatrix}$. We have

$$\min_{y \in \mathbb{R}^d} \|\hat{A}y - \hat{b}\|_2 \tag{10}$$

equivalent to $\|b - Ax\|_2^2 + \lambda \|Lx\|_2^2$, in the sense that for any $\hat{A}y \in \operatorname{range}(\hat{A})$, there is $x \in \mathbb{R}^d$ with $\hat{A}y = \begin{bmatrix} A \\ \sqrt{\lambda}L \end{bmatrix} x$, so that $\|\hat{A}y - \hat{b}\|_2^2 = \left\| \begin{bmatrix} A \\ \sqrt{\lambda}L \end{bmatrix} x - \hat{b} \right\|_2^2 = \|b - Ax\|_2^2 + \lambda \|Lx\|_2^2$. Let $y^* = \operatorname{argmin}_{y \in \mathbb{R}^d} \|\hat{A}y - \hat{b}\|_2$, so that $\hat{A}y^* = \begin{bmatrix} Ax^* \\ \sqrt{\lambda}Lx^* \end{bmatrix}$.

Let $\hat{A} = \begin{bmatrix} U_1 \\ U_2 \end{bmatrix}$, where $U_1 \in \mathbb{R}^{n \times d}$ and $U_2 \in \mathbb{R}^{d \times d}$, so that U_1 is as in the lemma statement.

We define \hat{S} to be $\begin{bmatrix} S & 0_{m \times d} \\ 0_{d \times n} & I_d \end{bmatrix}$ and \hat{S} satisfies Property (I) and (II) of Lemma 4.1.

Using $\|U_1^\top S^\top S U_1 - U_1^\top U_1\|_2 \leq 1/4$, with constant probability

$$\|\hat{A}^\top \hat{S}^\top \hat{S} \hat{A} - I_d\|_2 = \|U_1^\top S^\top S U_1 + U_2^\top U_2 - I_d\|_2 = \|U_1^\top S^\top S U_1 - U_1^\top U_1\|_2 \leq 1/4. \tag{11}$$

Using the normal equations for Eq. (10), we have

$$0 = \hat{A}^\top (\hat{b} - \hat{A}y^*) = U_1^\top (b - Ax^*) - \sqrt{\lambda} U_2^\top x^*,$$

and so

$$\hat{A}^\top \hat{S}^\top \hat{S} (\hat{b} - \hat{A}y^*) = U_1^\top S^\top S (b - Ax^*) - \sqrt{\lambda} U_2^\top x^* = U_1^\top S^\top S (b - Ax^*) - U_1^\top (b - Ax^*).$$

Using Property (II) of Lemma 4.1, with constant probability

$$\begin{aligned} & \|\hat{A}^\top \hat{S}^\top \hat{S} (\hat{b} - \hat{A}y^*)\|_2 \\ & = \|U_1^\top S^\top S (b - Ax^*) - U_1^\top (b - Ax^*)\|_2 \\ & \leq \sqrt{\epsilon \text{OPT}/2} \\ & = \sqrt{\epsilon/2} \|\hat{b} - \hat{A}y^*\|_2. \end{aligned} \tag{12}$$

It follows by a standard result from (11) and (12) that the solution $\tilde{y} \equiv \operatorname{argmin}_{y \in \mathbb{R}^d} \|\hat{S}(\hat{A}y - \hat{b})\|_2$ has $\|\hat{A}\tilde{y} - \hat{b}\|_2 \leq (1 + \epsilon) \min_{y \in \mathbb{R}^d} \|\hat{A}y - \hat{b}\|_2$, and therefore that \tilde{x} satisfies the claim of the theorem.

For convenience we give the proof of the standard result: (11) implies that $\hat{A}^\top \hat{S}^\top \hat{S} \hat{A}$ has smallest singular value at least $3/4$. The normal equations for the unsketched and sketched problems are

$$\hat{A}^\top (\hat{b} - \hat{A}y^*) = 0 = \hat{A}^\top \hat{S}^\top \hat{S} (\hat{b} - \hat{A}\tilde{y}).$$

The normal equations for the unsketched case imply $\|\hat{A}\tilde{y} - \hat{b}\|_2^2 = \|\hat{A}(\tilde{y} - y^*)\|_2^2 + \|\hat{b} - \hat{A}y^*\|_2^2$, so it is enough to show that $\|\hat{A}(\tilde{y} - y^*)\|_2^2 = \|\tilde{y} - y^*\|_2^2 \leq \epsilon \text{OPT}$. We have

$$\begin{aligned} (3/4)\|\tilde{y} - y^*\|_2 &\leq \|\hat{A}^\top \hat{S}^\top \hat{S} \hat{A}(\tilde{y} - y^*)\|_2 && \text{by Eq. (11)} \\ &= \|\hat{A}^\top \hat{S}^\top \hat{S} \hat{A}(\tilde{y} - y^*) - \hat{A}^\top \hat{S}^\top \hat{S} (\hat{b} - \hat{A}\tilde{y})\|_2 && \text{by Normal Equation} \\ &= \|\hat{A}^\top \hat{S}^\top \hat{S} (\hat{b} - \hat{A}y^*)\|_2 \\ &\leq \sqrt{\epsilon \text{OPT} / 2} && \text{by Eq. (12),} \end{aligned}$$

so that $\|\tilde{y} - y^*\|_2^2 \leq (4/3)^2 \epsilon \text{OPT} / 2 \leq \epsilon \text{OPT}$. The lemma follows. \square

The following lemma computes the statistical dimension $\text{sd}_\lambda(A, L)$ that will be used for computing the number of rows of sketching matrix S .

Lemma C.1. *For U_1 as in Lemma 4.1, $\|U_1\|_F^2 = \text{sd}_\lambda(A, L) = \sum_i 1/(1 + \lambda/\gamma_i^2) + d - p$, where A has singular values σ_i . Also $\|U_1\|_2 = \max\{1/\sqrt{1 + \lambda/\gamma_1^2}, 1\}$.*

Proof. Suppose we have the GSVD of (A, L) . Let

$$D \equiv \begin{bmatrix} \Sigma^\top \Sigma + \lambda \Omega^\top \Omega & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & I_{d-p} \end{bmatrix}^{-1/2}.$$

Then

$$\hat{A} = \begin{bmatrix} U \begin{bmatrix} \Sigma & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & I_{d-p} \end{bmatrix} D \\ \sqrt{\lambda} V \begin{bmatrix} \Omega & 0_{p \times (n-p)} \end{bmatrix} D \end{bmatrix}$$

has $\hat{A}^\top \hat{A} = I_d$, and for given x , there is $y = D^{-1} R Q^\top x$ with $\hat{A}y = \begin{bmatrix} A \\ \sqrt{\lambda} L \end{bmatrix} x$. We have $\|U_1\|_F^2 = \left\| U \begin{bmatrix} \Sigma & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & I_{d-p} \end{bmatrix} D \right\|_F^2 = \left\| \begin{bmatrix} \Sigma & 0_{p \times (n-p)} \\ 0_{(n-p) \times p} & I_{d-p} \end{bmatrix} D \right\|_F^2 = \sum_{i=1}^p 1/(1 + \lambda/\gamma_i^2) + d - p$ as claimed. \square

Theorem 4.3. (P-Spline regression) There is a constant $K > 0$ such that for $m \geq K(\epsilon^{-1} \text{sd}_\lambda(A, L) + \text{sd}_\lambda(A, L)^2)$ and $S \in \mathbb{R}^{m \times n}$ a sparse embedding matrix (e.g., COUNTSKETCH) with SA computable in $O(\text{nnz}(A))$ time, Property (I) and (II) of Lemma 4.1 apply, and with constant probability the corresponding $\tilde{x} = \operatorname{argmin}_{x \in \mathbb{R}^d} \|S(Ax - b)\|_2 + \lambda \|Lx\|_2^2$ is an ϵ -approximate solution to $\min_{x \in \mathbb{R}^d} \|b - Ax\|_2^2 + \lambda \|Lx\|_2^2$.

Proof. Recall that $\text{sd}_\lambda(A, L) = \|U_1\|_F^2$. Sparse embedding distributions satisfy the bound for approximate matrix multiplication

$$\|W^\top S^\top S H - W^\top H\|_F \leq C \|W\|_F \|H\|_F / \sqrt{m},$$

for a constant C (Clarkson and Woodruff, 2013; Meng and Mahoney, 2013; Nelson and Nguyen, 2013); this is also true of OSE matrices. We set $W = H = U_1$ and use $\|X\|_2 \leq \|X\|_F$ for all X and $m \geq K \|U_1\|_F^4$ to obtain Property (I) of Lemma 4.1, and set $W = U_1$, $H = b - Ax^*$ and use $m \geq K \|U_1\|_F^2 / \epsilon$ to obtain Property (II) of Lemma 4.1. (Here the bound is slightly stronger than Property (II), holding for $\lambda = 0$.) With Property (I) and Property (II), the claim for \tilde{x} from a sparse embedding follows using Lemma 4.1. \square

C.3 Proofs for Tensor Product ℓ_1 Regression

Lemma 5.3. *For any $p \geq 1$. Condition (\mathcal{A}) computes $\mathcal{AU}/(d\gamma_p)$ which is an $(\alpha, \beta\sqrt{3}d(tw)^{|1/p-1/2|}, p)$ -well-conditioned basis of \mathcal{A} , with probability at least $1 - \prod_{i=1}^q (n_i/w_i)\delta$.*

Proof. This lemma is similar to arguments in Clarkson et al. (2013), we simply adjust notation and parameters for completeness. Applying Theorem 5.2, we have that with probability at least $1 - \prod_{i=1}^q (n_i/w_i)\delta$, for all $x \in \mathbb{R}^r$, if we consider $y = \mathcal{A}x$ and write $y^\top = [z_1^\top, z_2^\top, \dots, z_{\prod_{i=1}^q n_i/w_i}^\top]^\top$, then for all $i \in [\prod_{i=1}^q n_i/w_i]$,

$$\sqrt{\frac{1}{2}}\|z_i\|_2 \leq \|S_i z_i\|_2 \leq \sqrt{\frac{3}{2}}\|z_i\|_2,$$

where $S_i \in \mathbb{R}^{m_i \times \prod_{i=1}^q w_i}$. In the following, suppose $m_i = t$. By relating the 2-norm and the p -norm, for $1 \leq p \leq 2$, we have

$$\|S_i z_i\|_p \leq t^{1/p-1/2} \|S_i z_i\|_2 \leq t^{1/p-1/2} \sqrt{\frac{3}{2}} \|z_i\|_2 \leq t^{1/p-1/2} \sqrt{\frac{3}{2}} \|z_i\|_p,$$

and similarly,

$$\|S_i z_i\|_p \geq \|S_i z_i\|_2 \geq \sqrt{\frac{1}{2}} \|z_i\|_2 \geq \sqrt{\frac{1}{2}} w^{1/2-1/p} \|z_i\|_p, \quad w = \prod_{j=1}^q w_j.$$

If $p > 2$, then

$$\|S_i z_i\|_p \leq \|S_i z_i\|_2 \leq \sqrt{\frac{3}{2}} \|z_i\|_2 \leq \sqrt{\frac{3}{2}} w^{1/2-1/p} \|z_i\|_p,$$

and similarly,

$$\|S_i z_i\|_p \geq t^{1/p-1/2} \|S_i z_i\|_2 \geq t^{1/p-1/2} \sqrt{\frac{1}{2}} \|z_i\|_2 \geq t^{1/p-1/2} \sqrt{\frac{1}{2}} \|z_i\|_p.$$

Since $\|\mathcal{A}x\|_p^p = \|y\|_p^p = \sum_i \|z_i\|_p^p$ and $\|S\mathcal{A}x\|_p^p = \sum_i \|S_i z_i\|_p^p$, for $p \in [1, 2]$ we have with probability $1 - \prod_{i=1}^q (n_i/w_i)\delta$

$$\sqrt{\frac{1}{2}} w^{1/2-1/p} \|\mathcal{A}x\|_p \leq \|S\mathcal{A}x\|_p \leq \sqrt{\frac{3}{2}} t^{1/p-1/2} \|\mathcal{A}x\|_p,$$

and for $p \in [2, \infty)$ with probability $1 - \prod_{i=1}^q (n_i/w_i)\delta$

$$\sqrt{\frac{1}{2}} t^{1/p-1/2} \|\mathcal{A}x\|_p \leq \|S\mathcal{A}x\|_p \leq \sqrt{\frac{3}{2}} w^{1/2-1/p} \|\mathcal{A}x\|_p.$$

In either case,

$$\|\mathcal{A}x\|_p \leq \gamma_p \|S\mathcal{A}x\|_p \leq \sqrt{3}(tw)^{|1/p-1/2|} \|\mathcal{A}x\|_p. \quad (13)$$

We have, from the definition of an (α, β, p) -well-conditioned basis, that

$$\|S\mathcal{A}U\|_p \leq \alpha \quad (14)$$

and for all $x \in \mathbb{R}^d$,

$$\|x\|_q \leq \beta \|S\mathcal{A}Ux\|_p. \quad (15)$$

Combining (13) and (14), we have that with probability at least $1 - \prod_{i=1}^q (n_i/w_i)\delta$,

$$\|\mathcal{AU}/(r\gamma_p)\|_p \leq \sum_i \|\mathcal{AU}_i/r\gamma_p\|_p \leq \sum_i \|S\mathcal{A}U_i/r\|_p \leq \alpha.$$

Combining (13) and (15), we have that with probability at least $1 - \prod_{i=1}^q (n_i/w_i)\delta$, for all $x \in \mathbb{R}^r$,

$$\|x\|_q \leq \beta \|S\mathcal{A}Ux\|_p \leq \beta \sqrt{3}r(tw)^{|1/p-1/2|} \|\mathcal{AU} \frac{1}{r\gamma_p} x\|_p.$$

Hence $\mathcal{AU}/(r\gamma_p)$ is an $(\alpha, \beta\sqrt{3}r(tw)^{|1/p-1/2|}, p)$ -well-conditioned basis. \square

Theorem 5.4. (Main result) Given $\epsilon \in (0, 1)$, $\mathcal{A} \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$, Alg. 3 computes \hat{x} such that with probability at least $1/2$, $\|\mathcal{A}\hat{x} - b\|_1 \leq (1 + \epsilon) \min_{x \in \mathbb{R}^d} \|\mathcal{A}x - b\|_1$. For the special case when $q = 2$, $n_1 = n_2$, the algorithm's running time is $O(n_1^{3/2} \text{poly}(\prod_{i=1}^2 d_i/\epsilon))$.

Proof. For notational simplicity, let us denote $n_{[q_1]} = \prod_{i=1}^{q_1} n_i$, $n_{[q] \setminus [q_1]} = \prod_{i=q_1+1}^q n_i$, $d_{[q_1]} = \prod_{i=1}^{q_1} d_i$, and $d_{[q] \setminus [q_1]} = \prod_{i=q_1+1}^q d_i$. For any row-block $A_{i_1}^1 \otimes \dots \otimes A_{i_q}^{(q)}$, computing $S_{i_1 i_2 \dots i_q} (A_{i_1}^1 \otimes \dots \otimes A_{i_q}^{(q-1)})$ takes $O(d(\sum_{k=1}^q \text{nnz}(A_{i_k}^{(k)})) + dqm \log(m))$ (see Sec 2). Hence for SA , it takes:

$$\left(d \sum_{k=1}^q \text{nnz}(A_k) \prod_{i \in [q] \setminus \{k\}} n_i/w_i \right) + \left(dqm \log(m) \prod_{i=1}^q n_i/w_i \right)$$

where $S \in \mathbb{R}^{(m \prod_{i=1}^q (n_i/w_i)) \times \prod_{i=1}^q w_i}$ and $m \geq 100 \prod_{i=1}^q d_i^2 (2 + 3^q)/\epsilon^2 = O(\text{poly}(d/\epsilon))$. We need to compute an orthogonal factorization $SA = QR_{\mathcal{A}}$ in $O(qmd^2)$ and then compute $U = R_{\mathcal{A}}^{-1}$ in $O(d^3)$ time. Hence the total running time of Algorithm Condition(A) is $O(qmd^2 + d^3)$. Thus the total running time of computing SA and Condition(A) is

$$O \left(\left(\sum_{k=1}^q \text{nnz}(A_k) \prod_{i \in [q] \setminus \{k\}} n_i/w_i \right) + \left(\prod_{i=1}^q n_i/w_i \right) \text{poly}(d/\epsilon) + qmd^2 + d^3 \right),$$

We will compute UG in $O(d^2 \log n)$ time. We compute $\tilde{E} = E(A_{q_1+1} \otimes \dots \otimes A_q)^T$ in $O(dn_{[q] \setminus [q_1]})$ time.

Then we can compute $R(A_1 \otimes \dots \otimes A_{q_1}) \tilde{E}_j$ in $O(n_{[q_1]} d_{[q_1]} \log n + d_{[q_1]} n_{[q] \setminus [q_1]} \log n)$ time.

Since computation of the median λ_i takes $O(\log n)$ time, computing all λ_i and then λ_e takes $O(n_{[q] \setminus [q_1]} \log n)$ time.

As AUG has $O(\log n)$ columns, we need to compute λ_e for each AUG using the above procedure and hence it takes in total $O(d(n_{[q_1]} + n_{[q] \setminus [q_1]}) \log^2 n)$ time.

Sampling a column of AUG using λ_e takes $O(\log n)$ time, sampling an entry in M takes in total $O(n_{[q_1]} + n_{[q] \setminus [q_1]})$ time.

Since we need $\sqrt{\prod_{k=1}^q w_k} \text{poly}(r)$ samples to select rows, the running time is $d(n_{[q_1]} + n_{[q] \setminus [q_1]}) \log^2 n \cdot \sqrt{\prod_{k=1}^q w_k} \text{poly}(r)$.

Now for simplicity, we set $q = 2$, $n_i = n_0$ for $i \in [2]$. Note that it is optimal to choose $w_i = w$ for $i \in [2]$. Substituting $q = 2, n_i = n_0$ and $w_i = w$, we that the total running time of Alg. 3:

$$O(dw^{-1}n_0(\text{nnz}(A_1) + \text{nnz}(A_2)) + w^{-2}n_0^2 \text{poly}(d/\epsilon) + wn_0 \text{poly}(d) \log(n)).$$

For dense A_1 and A_2 , $\text{nnz}(A_1) + \text{nnz}(A_2) = O(n_0)$ time, and so ignoring poly and \log terms that do not depend on n_0 , the total running time can be simplified to:

$$O(w^{-1}n_0^2 + wn_0).$$

Setting $w = \sqrt{n_0}$, we can minimize the above running time to $O(n_0^{3/2})$, which is faster than the n_0^2 time for solving the problem by forming $A_1 \otimes A_2$. \square