# Fast Regression with an $\ell_\infty$ Guarantee[*]

Eric Price
zhaos@utexas.edu
UT-Austin

Zhao Song
zhaos@utexas.edu
UT-Austin

David P. Woodruff
dpwoodru@us.ibm.com
IBM Almaden

May 31, 2017

## Abstract

Sketching has emerged as a powerful technique for speeding up problems in numerical linear algebra, such as regression. In the overconstrained regression problem, one is given an $n \times d$ matrix $A$, with $n \gg d$, as well as an $n \times 1$ vector $b$, and one wants to find a vector $\widehat{x}$ so as to minimize the residual error $\|Ax - b\|_2$. Using the sketch and solve paradigm, one first computes $S \cdot A$ and $S \cdot b$ for a randomly chosen matrix $S$, then outputs $x' = (SA)^\dagger Sb$ so as to minimize $\|SAx' - Sb\|_2$.

The sketch-and-solve paradigm gives a bound on $\|x' - x^*\|_2$ when $A$ is well-conditioned. Our main result is that, when $S$ is the subsampled randomized Fourier/Hadamard transform, the error $x' - x^*$ behaves as if it lies in a "random" direction within this bound: for any fixed direction $a \in \mathbb{R}^d$, we have with $1 - d^{-c}$ probability that

$$\langle a, x' - x^* \rangle \lesssim \frac{\|a\|_2 \|x' - x^*\|_2}{d^{\frac{1}{2} - \gamma}}, \tag{1}$$

where $c, \gamma > 0$ are arbitrary constants. This implies $\|x' - x^*\|_\infty$ is a factor $d^{\frac{1}{2} - \gamma}$ smaller than $\|x' - x^*\|_2$. It also gives a better bound on the generalization of $x'$ to new examples: if rows of $A$ correspond to examples and columns to features, then our result gives a better bound for the error introduced by sketch-and-solve when classifying fresh examples. We show that not all oblivious subspace embeddings $S$ satisfy these properties. In particular, we give counterexamples showing that matrices based on Count-Sketch or leverage score sampling do not satisfy these properties.

We also provide lower bounds, both on how small $\|x' - x^*\|_2$ can be, and for our new guarantee (1), showing that the subsampled randomized Fourier/Hadamard transform is nearly optimal. Our lower bound on $\|x' - x^*\|_2$ shows that there is an $O(1/\varepsilon)$ separation in the dimension of the optimal oblivious subspace embedding required for outputting an $x'$ for which $\|x' - x^*\|_2 \le \epsilon \|Ax^* - b\|_2 \cdot \|A^\dagger\|_2$, compared to the dimension of the optimal oblivious subspace embedding required for outputting an $x'$ for which $\|Ax' - b\|_2 \le (1 + \epsilon)\|Ax^* - b\|_2$, that is, the former problem requires dimension $\Omega(d/\epsilon^2)$ while the latter problem can be solved with dimension $O(d/\epsilon)$. This explains the reason known upper bounds on the dimensions of these two variants of regression have differed in prior work.

---

# 1   Introduction

Oblivious subspace embeddings (OSEs) were introduced by Sarlos [Sar06] to solve linear algebra problems more quickly than traditional methods. An OSE is a distribution of matrices $S \in \mathbb{R}^{m \times n}$ with $m \ll n$ such that, for any $d$-dimensional subspace $U \subset \mathbb{R}^n$, with "high" probability $S$ preserves the norm of every vector in the subspace. OSEs are a generalization of the classic Johnson-Lindenstrauss lemma from vectors to subspaces. Formally, we require that with probability $1 - \delta$,

$$\|Sx\|_2 = (1 \pm \varepsilon)\|x\|_2$$

simultaneously for all $x \in U$, that is, $(1 - \varepsilon)\|x\|_2 \le \|Sx\|_2 \le (1 + \varepsilon)\|x\|_2$.

A major application of OSEs is to regression. The regression problem is, given $b \in \mathbb{R}^n$ and $A \in \mathbb{R}^{n \times d}$ for $n \ge d$, to solve for

$$x^* = \arg\min_{x \in \mathbb{R}^d} \|Ax - b\|_2 \tag{2}$$

Because $A$ is a "tall" matrix with more rows than columns, the system is overdetermined and there is likely no solution to $Ax = b$, but regression will find the closest point to $b$ in the space spanned by $A$. The classic answer to regression is to use the Moore-Penrose pseudoinverse: $x^* = A^\dagger b$ where

$$A^\dagger = (A^\top A)^{-1} A^\top$$

is the "pseudoinverse" of $A$ (assuming $A$ has full column rank, which we will typically do for simplicity). This classic solution takes $O(nd^{\omega-1} + d^\omega)$ time, where $\omega < 2.373$ is the matrix multiplication constant [CW90, Wil12, Gal14]: $nd^{\omega-1}$ time to compute $A^\top A$ and $d^\omega$ time to compute the inverse.

OSEs speed up the process by replacing (2) with

$$x' = \arg\min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2$$

for an OSE $S$ on $d+1$-dimensional spaces. This replaces the $n \times d$ regression problem with an $m \times d$ problem, which can be solved more quickly since $m \ll n$. Because $Ax - b$ lies in the $d + 1$-dimensional space spanned by $b$ and the columns of $A$, with high probability $S$ preserves the norm of $SAx - Sb$ to $1 \pm \varepsilon$ for all $x$. Thus,

$$\|Ax' - b\|_2 \le \frac{1+\varepsilon}{1-\varepsilon}\|Ax^* - b\|_2.$$

That is, $S$ produces a solution $x'$ which preserves the *cost* of the regression problem. The running time for this method depends on (1) the reduced dimension $m$ and (2) the time it takes to multiply $S$ by $A$. We can compute these for "standard" OSE types:

- If $S$ has i.i.d. Gaussian entries, then $m = O(d/\varepsilon^2)$ is sufficient (and in fact, $m \ge d/\epsilon^2$ is required [NN14]). However, computing $SA$ takes $O(mnd) = O(nd^2/\varepsilon^2)$ time, which is worse than solving the original regression problem (one can speed this up using fast matrix multiplication, though it is still worse than solving the original problem).

- If $S$ is a subsampled randomized Hadamard transform (SRHT) matrix with random sign flips (see Theorem 2.4 in [Woo14] for a survey, and also see [CNW16] which gives a recent improvement) then $m$ increases to $\widetilde{O}(d/\varepsilon^2 \cdot \log n)$, where $\widetilde{O}(f) = f\text{poly}(\log(f))$. But now, we can compute $SA$ using the fast Hadamard transform in $O(nd \log n)$ time. This makes the overall regression problem take $O(nd \log n + d^\omega/\varepsilon^2)$ time.

- If $S$ is a random sparse matrix with random signs (the "Count-Sketch" matrix), then $m = d^{1+\gamma}/\varepsilon^2$ suffices for $\gamma > 0$ a decreasing function of the sparsity [CW13, MM13, NN13, BDN15, Coh16]. (The definition of a Count-Sketch matrix is, for any $s \geq 1$, $S_{i,j} \in \{0, -1/\sqrt{s}, 1/\sqrt{s}\}$, $\forall i \in [m], j \in [n]$ and the column sparsity of matrix $S$ is $s$. Independently in each column $s$ positions are chosen uniformly at random without replacement, and each chosen position is set to $-1/\sqrt{s}$ with probability $1/2$, and $+1/\sqrt{s}$ with probability $1/2$.) Sparse OSEs can benefit from the sparsity of $A$, allowing for a running time of $\widetilde{O}(\text{nnz}(A)) + \widetilde{O}(d^\omega/\varepsilon^2)$, where $\text{nnz}(A)$ denotes the number of non-zeros in $A$.

When $n$ is large, the latter two algorithms are substantially faster than the naïve $nd^{\omega-1}$ method.

## 1.1 Our Contributions

Despite the success of using subspace embeddings to speed up regression, often what practitioners are interested is not in preserving the cost of the regression problem, but rather in the *generalization* or *prediction* error provided by the vector $x'$. Ideally, we would like for any future (unseen) example $a \in \mathbb{R}^d$, that $\langle a, x' \rangle \approx \langle a, x^* \rangle$ with high probability.

Ultimately one may want to use $x'$ to do classification, such as regularized least squares classification (RLSC) [RYP03], which has been found in cases to do as well as support vector machines but is much simpler [ZP04]. In this application, given a training set of examples with multiple (non-binary) labels identified with the rows of an $n \times d$ matrix $A$, one creates an $n \times r$ matrix $B$, each column indicating the presence or absence of one of the $r$ possible labels in each example. One then solves the multiple response regression problem $\min_X \|AX - B\|_F$, and uses $X$ to classify future examples. A commonly used method is for a future example $a$, to compute $\langle a, x_1 \rangle, \ldots, \langle a, x_r \rangle$, where $x_1, \ldots, x_r$ are the columns of $X$. One then chooses the label $i$ for which $\langle a, x_i \rangle$ is maximum.

For this to work, we would like the inner products $\langle a, x'_1 \rangle, \ldots, \langle a, x'_r \rangle$ to be close to $\langle a, x^*_1 \rangle, \ldots, \langle a, x^*_r \rangle$, where $X'$ is the solution to $\min_X \|SAX - SB\|_F$ and $X^*$ is the solution to $\min_X \|AX - B\|_F$. For any $O(1)$-accurate OSE on $d + r$ dimensional spaces [Sar06], which also satisfies so-called approximate matrix multiplication with error $\varepsilon' = \varepsilon/\sqrt{(d+r)}$, we get that

$$\|x' - x^*\|_2 \leq O(\varepsilon) \cdot \|Ax^* - b\|_2 \cdot \|A^\dagger\|_2 \tag{3}$$

where $\|A^\dagger\|$ is the spectral norm of $A^\dagger$, which equals the reciprocal of the smallest singular value of $A$. To obtain a generalization error bound for an unseen example $a$, one has

$$|\langle a, x^* \rangle - \langle a, x' \rangle| = |\langle a, x^* - x' \rangle| \leq \|x^* - x'\|_2 \|a\|_2 = O(\varepsilon)\|a\|_2 \|Ax^* - b\|_2 \|A^\dagger\|_2, \tag{4}$$

which could be tight if given only the guarantee in (3). However, if the difference vector $x' - x^*$ were distributed in a uniformly random direction subject to (3), then one would expect an $\widetilde{O}(\sqrt{d})$ factor improvement in the bound. This is what our main theorem shows:

**Theorem 1** (Main Theorem, informal). *Suppose $n \leq \text{poly}(d)$ and matrix $A \in \mathbb{R}^{n \times d}$ and vector $b \in \mathbb{R}^n$ are given. Let $S \in \mathbb{R}^{m \times n}$ be a subsampled randomized Hadamard transform matrix with $m = d^{1+\gamma}/\varepsilon^2$ rows for an arbitrarily small constant $\gamma > 0$. For $x' = \arg\min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2$ and $x^* = \arg\min_{x \in \mathbb{R}^d} \|Ax - b\|_2$, and any fixed $a \in \mathbb{R}^d$,*

$$|\langle a, x^* \rangle - \langle a, x' \rangle| \leq \frac{\varepsilon}{\sqrt{d}} \|a\|_2 \|Ax^* - b\|_2 \|A^\dagger\|_2. \tag{5}$$

*with probability $1 - 1/d^C$ for an arbitrarily large constant $C > 0$. This implies that*

$$\|x^* - x'\|_\infty \leq \frac{\varepsilon}{\sqrt{d}} \|Ax^* - b\|_2 \|A^\dagger\|_2. \tag{6}$$

*with $1 - 1/d^{C-1}$ probability.*

*If $n > \text{poly}(d)$, then by first composing $S$ with a Count-Sketch OSE with $\text{poly}(d)$ rows, one can achieve the same guarantee.*

(Here $\gamma$ is a constant going to zero as $n$ increases; see Theorem 10 for a formal statement of Theorem 1.)

Notice that Theorem 1 is considerably stronger than that of (4) provided by existing guarantees. Indeed, in order to achieve the guarantee (6) in Theorem 1, one would need to set $\varepsilon' = \varepsilon/\sqrt{d}$ in existing OSEs, resulting in $\Omega(d^2/\epsilon^2)$ rows. In contrast, we achieve only $d^{1+\gamma}/\epsilon^2$ rows. We can improve the bound in Theorem 1 to $m = d/\varepsilon^2$ if $S$ is a matrix of i.i.d. Gaussians; however, as noted, computing $S \cdot A$ is slower in this case.

Note that Theorem 1 also *makes no distributional assumptions* on the data, and thus the data could be heavy-tailed or even adversarially corrupted. This implies that our bound is still useful when the rows of $A$ are not sampled independently from a distribution with bounded variance.

The $\ell_\infty$ bound (6) of Theorem 1 is achieved by applying (5) to the standard basis vectors $a = e_i$ for each $i \in [d]$ and applying a union bound. This $\ell_\infty$ guarantee often has a more natural interpretation than the $\ell_2$ guarantee—if we think of the regression as attributing the observable as a sum of various factors, (6) says that the contribution of each factor is estimated well. One may also see our contribution as giving a way for estimating the pseudoinverse $A^\dagger$ *entrywise*. Namely, we get that $(SA)^\dagger S \approx A^\dagger$ in the sense that each entry is within additive $O(\varepsilon\sqrt{\frac{\log d}{d}}\|A^\dagger\|_2)$. There is a lot of work on computing entries of inverses of a matrix, see, e.g., [ADL$^+$12, LAKD08].

Another benefit of the $\ell_\infty$ guarantee is when the regression vector $x^*$ is expected to be $k$-*sparse* (e.g. [Lee12]). In such cases, thresholding to the top $k$ entries will yield an $\ell_2$ guarantee a factor $\sqrt{\frac{k}{d}}$ better than (3).

One could ask if Theorem 1 also holds for sparse OSEs, such as the Count-Sketch. Surprisingly, we show that one cannot achieve the generalization error guarantee in Theorem 1 with high probability, say, $1 - 1/d$, using such embeddings, despite the fact that such embeddings do approximate the cost of the regression problem up to a $1 + \epsilon$ factor with high probability. This shows that the generalization error guarantee is achieved by some subspace embeddings but not all.

**Theorem 2** (Not all subspace embeddings give the $\ell_\infty$ guarantee; informal version of Theorem 20)**.** *The Count-Sketch matrix with $d^{1.5}$ rows and sparsity $d^{.25}$—which is an OSE with exponentially small failure probability—with constant probability will have a result $x'$ that does not satisfy the $\ell_\infty$ guarantee (6).*

We can show that Theorem 1 holds for $S$ based on the Count-Sketch OSE $T$ with $d^{O(C)}/\epsilon^2$ rows with $1 - 1/d^C$ probability. We can thus compose the Count-Sketch OSE with the SRHT matrix and obtain an $O(\text{nnz}(A)) + \text{poly}(d/\epsilon)$ time algorithm to compute $S \cdot TA$ achieving (6). We can also compute $R \cdot S \cdot T \cdot A$, where $R$ is a matrix of Gaussians, which is more efficient now that $STA$ only has $d^{1+\gamma}/\epsilon^2$ rows; this will reduce the number of rows to $d/\epsilon^2$.

Another common method of dimensionality reduction for linear regression is *leverage score sampling* [DMIMW12, LMP13, PKB14, CMM15], which subsamples the rows of $A$ by choosing each row with probability proportional to its "leverage scores". With $O(d\log(d/\delta)/\varepsilon^2)$ rows taken, the result $x'$ will satisfy the $\ell_2$ bound (3) with probability $1 - \delta$. However, it does not give a good $\ell_\infty$ bound:

**Theorem 3** (Leverage score sampling does not give the $\ell_\infty$ guarantee; informal version of Theorem 23)**.** *Leverage score sampling with $d^{1.5}$ rows—which satisfies the $\ell_2$ bound with exponentially*

*small failure probability—with constant probability will have a result $x'$ that does not satisfy the $\ell_\infty$ guarantee ([6]).*

Finally, we show that the $d^{1+\gamma}/\varepsilon^2$ rows that SRHT matrices use is roughly optimal:

**Theorem 4** (Lower bounds for $\ell_2$ and $\ell_\infty$ guarantees; informal versions of of Theorem [14] and Corollary [18]). *Any sketching matrix distribution over $m \times n$ matrices that satisfies either the $\ell_2$ guarantee ([3]) or the $\ell_\infty$ guarantee ([6]) must have $m \gtrsim \min(n, d/\varepsilon^2)$.*

Notice that our result shows the necessity of the $1/\varepsilon$ separation between the results originally defined in Equation (3) and (4) of Theorem 12 of [Sar06]. If we want to output some vector $x'$ such that $\|Ax' - b\|_2 \leq (1+\varepsilon)\|Ax^* - b\|_2$, then it is known that $m = \Theta(d/\varepsilon)$ is necessary and sufficient. However, if we want to output a vector $x'$ such that $\|x' - x^*\|_2 \leq \varepsilon\|Ax^* - b\|_2 \cdot \|A^\dagger\|_2$, then we show that $m = \Theta(d/\varepsilon^2)$ is necessary and sufficient.

### 1.1.1 Comparison to Gradient Descent

While this work is primarily about sketching methods, one could instead apply iterative methods such as gradient descent, after appropriately preconditioning the matrix, see, e.g., [AMT10, ZF13, CW13]. That is, one can use an OSE with constant $\varepsilon$ to construct a preconditioner for $A$ and then run conjugate gradient using the preconditioner. This gives an overall dependence of $\log(1/\epsilon)$.

The main drawback of this approach is that one loses the ability to save on storage space or number of passes when $A$ appears in a stream, or to save on communication or rounds when $A$ is distributed. Given increasingly large data sets, such scenarios are now quite common, see, e.g., [CW09] for regression algorithms in the data stream model. In situations where the entries of $A$ appear sequentially, for example, a row at a time, one does not need to store the full $n \times d$ matrix $A$ but only the $m \times d$ matrix $SA$.

Also, iterative methods can be less efficient when solving multiple response regression, where one wants to minimize $\|AX - B\|$ for a $d \times t$ matrix $X$ and an $n \times t$ matrix $B$. This is the case when $\varepsilon$ is constant and $t$ is large, which can occur in some applications (though there are also other applications for which $\varepsilon$ is very small). For example, conjugate gradient with a preconditioner will take $\widetilde{O}(ndt)$ time while using an OSE directly will take only $\widetilde{O}(nd + d^2 t)$ time (since one effectively replaces $n$ with $O(d)$ after computing $S \cdot A$), separating $t$ from $d$. Multiple response regression, arises, for example, in the RLSC application above.

### 1.1.2 Proof Techniques

**Theorem [1].** As noted in Theorem [2], there are some OSEs for which our generalization error bound does not hold. This hints that our analysis is non-standard and cannot use generic properties of OSEs as a black box. Indeed, in our analysis, we have to consider matrix products of the form $S^\top S(UU^\top S^\top S)^k$ for our random sketching matrix $S$ and a fixed matrix $U$, where $k$ is a positive integer. We stress that it is the *same matrix $S$* appearing multiple times in this expression, which considerably complicates the analysis, and does not allow us to appeal to standard results on approximate matrix product (see, e.g., [Woo14] for a survey). The key idea is to recursively reduce $S^\top S(UU^\top S^\top S)^k$ using a property of $S$. We use properties that only hold for specifics OSEs $S$: first, that each column of $S$ is unit vector; and second, that for all pairs $(i,j)$ and $i \neq j$, the inner product between $S_i$ and $S_j$ is at most $\frac{\sqrt{\log n}}{\sqrt{m}}$ with probability $1 - 1/\text{poly}(n)$.

**Theorems [20] and [23].** To show that Count-Sketch does not give the $\ell_\infty$ guarantee, we construct a matrix $A$ and vector $b$ as in Figure [1], which has optimal solution $x^*$ with all coordinates

$1/\sqrt{d}$. We then show, for our setting of parameters, that there likely exists an index $j \in [d]$ satisfying the following property: the $j$th column of $S$ has disjoint support from the $k$th column of $S$ for all $k \in [d + \alpha] \setminus \{j\}$ except for a single $k > d$, for which $S_j$ and $S_k$ share exactly one common entry in their support. In such cases we can compute $x'_j$ explicitly, getting $|x'_j - x^*_j| = \frac{1}{s\sqrt{\alpha}}$. By choosing suitable parameters in our construction, this gives that $\|x' - x^*\|_\infty \gg \frac{1}{\sqrt{d}}$. The lower bound for leverage score sampling follows a similar construction.

**Theorem 14 and Corollary 18.** The lower bound proof for the $\ell_2$ guarantee uses Yao's minimax principle. We are allowed to fix an $m \times n$ sketching matrix $S$ and design a distribution over $[A\ b]$. We first write the sketching matrix $S = U\Sigma V^\top$ in its singular value decomposition (SVD). We choose the $d + 1$ columns of the adjoined matrix $[A, b]$ to be random orthonormal vectors. Consider an $n \times n$ orthonormal matrix $R$ which contains the columns of $V$ as its first $m$ columns, and is completed on its remaining $n - m$ columns to an arbitrary orthonormal basis. Then $S \cdot [A, b] = V^\top RR^\top \cdot [A, b] = [U\Sigma I_m, 0] \cdot [R^\top A, R^\top b]$. Notice that $[R^\top A, R^\top b]$ is equal in distribution to $[A, b]$, since $R$ is fixed and $[A, b]$ is a random matrix with $d + 1$ orthonormal columns. Therefore, $S \cdot [A, b]$ is equal in distribution to $[U\Sigma G, U\Sigma h]$ where $[G, h]$ corresponds to the first $m$ rows of an $n \times (d + 1)$ uniformly random matrix with orthonormal columns.

A key idea is that if $n = \Omega(\max(m, d)^2)$, then by a result of Jiang [J+06], any $m \times (d + 1)$ submatrix of a random $n \times n$ orthonormal matrix has $o(1)$ total variation distance to a $d \times d$ matrix of i.i.d. $N(0, 1/n)$ random variables, and so any events that would have occurred had $G$ and $h$ been independent i.i.d. Gaussians, occur with the same probability for our distribution up to an $1 - o(1)$ factor, so we can assume $G$ and $h$ are independent i.i.d. Gaussians in the analysis.

The optimal solution $x'$ in the sketch space equals $(SA)^\dagger Sb$, and by using that $SA$ has the form $U\Sigma G$, one can manipulate $\|(SA)^\dagger Sb\|$ to be of the form $\|\tilde{\Sigma}^\dagger (\Sigma R)^\dagger \Sigma h\|_2$, where the SVD of $G$ is $R\tilde{\Sigma}T$. We can upper bound $\|\tilde{\Sigma}\|_2$ by $\sqrt{r/n}$, since it is just the maximum singular value of a Gaussian matrix, where $r$ is the rank of $S$, which allows us to lower bound $\|\tilde{\Sigma}^\dagger (\Sigma R)^\dagger \Sigma h\|_2$ by $\sqrt{n/r}\|(\Sigma R)^\dagger \Sigma h\|_2$. Then, since $h$ is i.i.d. Gaussian, this quantity concentrates to $\frac{1}{\sqrt{r}}\|(\Sigma R)^\dagger \Sigma h\|$, since $\|Ch\|^2 \approx \|C\|_F^2/n$ for a vector $h$ of i.i.d. $N(0, 1/n)$ random variables. Finally, we can lower bound $\|(\Sigma R)^\dagger \Sigma\|_F^2$ by $\|(\Sigma R)^\dagger \Sigma RR^\top\|_F^2$ by the Pythagorean theorem, and now we have that $(\Sigma R)^\dagger \Sigma R$ is the identity, and so this expression is just equal to the rank of $\Sigma R$, which we prove is at least $d$. Noting that $x^* = 0$ for our instance, putting these bounds together gives $\|x' - x^*\| \geq \sqrt{d/r}$. The last ingredient is a way to ensure that the rank of $S$ is at least $d$. Here we choose another distribution on inputs $A$ and $b$ for which it is trivial to show the rank of $S$ is at least $d$ with large probability. We require $S$ be good on the mixture. Since $S$ is fixed and good on the mixture, it is good for both distributions individually, which implies we can assume $S$ has rank $d$ in our analysis of the first distribution above.

## 1.2 Notation

For a positive integer, let $[n] = \{1, 2, \ldots, n\}$. For a vector $x \in \mathbb{R}^n$, define $\|x\|_2 = (\sum_{i=1}^n x_i^2)^{\frac{1}{2}}$ and $\|x\|_\infty = \max_{i \in [n]} |x_i|$. For a matrix $A \in \mathbb{R}^{m \times n}$, define $\|A\|_2 = \sup_x \|Ax\|_2/\|x\|_2$ to be the spectral norm of $A$ and $\|A\|_F = (\sum_{i,j} A_{i,j}^2)^{1/2}$ to be the Frobenius norm of $A$. We use $A^\dagger$ to denote the Moore-Penrose pseudoinverse of $m \times n$ matrix $A$, which if $A = U\Sigma V^\top$ is its SVD (where $U \in \mathbb{R}^{m \times n}$, $\Sigma \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{n \times}$ for $m \geq n$), is given by $A^\dagger = V\Sigma^{-1}U^\top$.

In addition to $O(\cdot)$ notation, for two functions $f, g$, we use the shorthand $f \lesssim g$ (resp. $\gtrsim$) to indicate that $f \leq Cg$ (resp. $\geq$) for an absolute constant $C$. We use $f \approx g$ to mean $cf \leq g \leq Cf$ for constants $c, C$.
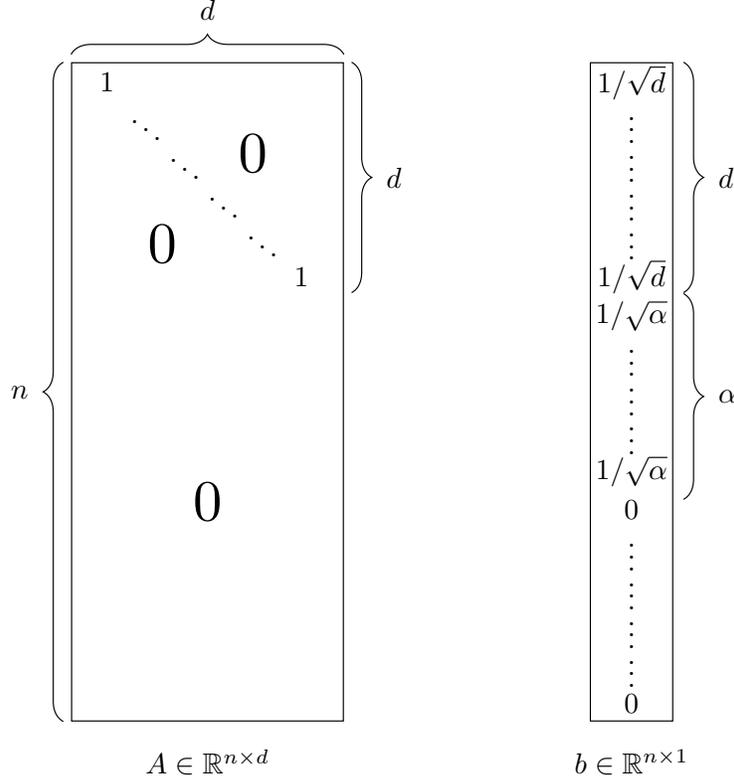
Figure 1: Our construction of $A$ and $b$ for the proof that Count-Sketch does not obey the $\ell_\infty$ guarantee. $\alpha < d$.

**Definition 5** (Subspace Embedding). *A $(1 \pm \epsilon)$ $\ell_2$-subspace embedding for the column space of an $n \times d$ matrix $A$ is a matrix $S$ for which for all $x \in \mathbb{R}^d$, $\|SAx\|_2^2 = (1 \pm \epsilon)\|Ax\|_2^2$.*

**Definition 6** (Approximate Matrix Product). *Let $0 < \epsilon < 1$ be a given approximation parameter. Given matrices $A$ and $B$, where $A$ and $B$ each have $n$ rows, the goal is to output a matrix $C$ so that $\|A^\top B - C\|_F \le \epsilon\|A\|_F\|B\|_F$. Typically $C$ has the form $A^\top S^\top SB$, for a random matrix $S$ with a small number of rows. In particular, this guarantee holds for the subsampled randomized Hadamard transform $S$ with $O(\epsilon^{-2})$ rows [DMMS11].*

## 2   Warmup: Gaussians OSEs

We first show that if $S$ is a Gaussian random matrix, then it satisfies the generalization guarantee. This follows from the rotational invariance of the Gaussian distribution.

**Theorem 7.** *Suppose $A \in \mathbb{R}^{n \times d}$ has full column rank. If the entries of $S \in \mathbb{R}^{m \times n}$ are i.i.d. $N(0, 1/m)$, $m = O(d/\varepsilon^2)$, then for any vectors $a, b$ and $x^* = A^\dagger b$, we have, with probability $1 - 1/\mathrm{poly}(d)$,*

$$|a^\top (SA)^\dagger Sb - a^\top x^*| \lesssim \frac{\varepsilon\sqrt{\log d}}{\sqrt{d}}\|a\|_2\|b - Ax^*\|_2\|A^\dagger\|_2.$$

Because $SA$ has full column rank with probability 1, $(SA)^\dagger SA = I$. Therefore

$$|a^\top (SA)^\dagger Sb - a^\top x^*| = |a^\top (SA)^\dagger S(b - Ax^*)| = |a^\top (SA)^\dagger S(b - AA^\dagger b)|.$$

Thus it suffices to only consider vectors $b$ where $A^\dagger b = 0$, or equivalently $U^\top b = 0$. In such cases, $SU$ will be independent of $Sb$, which will give the result. The proof is in Appendix A.

## 3 SRHT Matrices

We first provide the definition of the subsampled randomized Hadamard transform(SRHT): let $S = \frac{1}{\sqrt{rn}} P H_n D$. Here, $D$ is an $n \times n$ diagonal matrix with i.i.d. diagonal entries $D_{i,i}$, for which $D_{i,i}$ in uniform on $\{-1, +1\}$. The matrix $H_n$ is the Hadamard matrix of size $n \times n$, and we assume $n$ is a power of 2. Here, $H_n = [H_{n/2}, \ H_{n/2}; H_{n/2}, \ -H_{n/2}]$ and $H_1 = [1]$. The $r \times n$ matrix $P$ samples $r$ coordinates of an $n$ dimensional vector uniformly at random.

For other subspace embeddings, we no longer have that $SU$ and $Sb$ are independent. To analyze them, we start with a claim that allows us to relate the inverse of a matrix to a power series.

**Claim 8.** *Let $S \in \mathbb{R}^{m \times n}$, $A \in \mathbb{R}^{n \times d}$ have SVD $A = U\Sigma V^\top$, and define $T \in \mathbb{R}^{d \times d}$ by*

$$T = I_d - U^\top S^\top S U.$$

*Suppose $SA$ has linearly independent columns and $\|T\|_2 \le 1/2$. Then*

$$(SA)^\dagger S = V\Sigma^{-1} \left( \sum_{k=0}^{\infty} T^k \right) U^\top S^\top S. \tag{7}$$

*Proof.*

$$\begin{aligned}
(SA)^\dagger S &= (A^\top S^\top S A)^{-1} A^\top S^\top S \\
&= (V\Sigma U^\top S^\top S U \Sigma V^\top)^{-1} V\Sigma U^\top S^\top S \\
&= V\Sigma^{-1} (U^\top S^\top S U)^{-1} U^\top S^\top S \\
&= V\Sigma^{-1} (I_d - T)^{-1} U^\top S^\top S \\
&= V\Sigma^{-1} \left( \sum_{k=0}^{\infty} T^k \right) U^\top S^\top S,
\end{aligned}$$

where in the last equality, since $\|T\|_2 < 1$, the von Neumann series $\sum_{k=0}^{\infty} T^k$ converges to $(I_d - T)^{-1}$. ∎

We then bound the $k$th term of this sum:

**Lemma 9.** *Let $S \in \mathbb{R}^{r \times n}$ be the subsampled randomized Hadamard transform, and let $a$ be a unit vector. Then with probability $1 - 1/\mathrm{poly}(n)$, we have*

$$|a^\top S^\top S (U U^\top S^\top S)^k b| = O(\log^k n) \cdot (O(d(\log n)/r) + 1)^{\frac{k-1}{2}} \cdot (\sqrt{d}\|b\|_2 (\log n)/r + \|b\|_2 (\log^{\frac{1}{2}} n)/r^{\frac{1}{2}})$$

*Hence, for $r$ at least $d \log^{2k+2} n \log^2(n/\varepsilon)/\varepsilon^2$, this is at most $O(\|b\|_2 \varepsilon/\sqrt{d})$ with probability at least $1 - 1/\mathrm{poly}(n)$.*

We defer the proof of this lemma to the next section, and now show how the lemma lets us prove that SRHT matrices satisfy the generalization bound with high probability:

7

**Theorem 10.** *Suppose $A \in \mathbb{R}^{n \times d}$ has full column rank with $\log n = d^{o(1)}$. Let $S \in \mathbb{R}^{m \times n}$ be a subsampled randomized Hadamard transform with $m = O(d^{1+\alpha}/\varepsilon^2)$ for $\alpha = \Theta(\sqrt{\frac{\log\log n}{\log d}})$. For any vectors $a, b$ and $x^* = A^\dagger b$, we have*

$$|a^\top (SA)^\dagger S b - a^\top x^*| \lesssim \frac{\varepsilon}{\sqrt{d}} \|a\|_2 \|b - A x^*\|_2 \|\Sigma^{-1}\|_2$$

*with probability $1 - 1/\text{poly}(d)$.*

*Proof.* Define $\Delta = \Theta\left(\frac{1}{\sqrt{m}}\right)(\log^c d)\|a\|_2 \|b - A x^*\|_2 \|\Sigma^{-1}\|_2$. For a constant $c > 0$, we have that $S$ is a $(1 \pm \gamma)$-subspace embedding (Definition 5) for $\gamma = \sqrt{\frac{d \log^c n}{m}}$ with probability $1 - 1/\text{poly}(d)$ (see, e.g., Theorem 2.4 of [Woo14] and references therein), so $\|SUx\|_2 = (1 \pm \gamma)\|Ux\|_2$ for all $x$, which we condition on. Hence for $T = I_d - U^\top S^\top S U$, we have $\|T\|_2 \leq (1 + \gamma)^2 - 1 \lesssim \gamma$. In particular, $\|T\|_2 < 1/2$ and we can apply Claim 8.

As in Section 2, $SA$ has full column rank if $S$ is a subspace embedding, so $(SA)^\dagger SA = I$ and we may assume $x^* = 0$ without loss of generality.

By the approximate matrix product (Definition 6), we have for some $c$ that

$$|a^\top V \Sigma^{-1} U^\top S^\top S b| \leq \frac{\log^c d}{\sqrt{m}} \|a\|_2 \|b\|_2 \|\Sigma^{-1}\|_2 \leq \Delta \tag{8}$$

with $1 - 1/\text{poly}(d)$ probability. Suppose this event occurs, bounding the $k = 0$ term of (7). Hence it suffices to show that the $k \geq 1$ terms of (7) are bounded by $\Delta$.

By approximate matrix product (Definition 6), we also have with $1 - 1/d^2$ probability that

$$\|U^\top S^\top S b\|_F \leq \frac{\log^c d}{\sqrt{m}} \|U^\top\|_F \|b\|_2 \leq \frac{\log^c d \sqrt{d}}{\sqrt{m}} \|b\|_2.$$

Combining with $\|T\|_2 \lesssim \gamma$ we have for any $k$ that

$$|a^\top V \Sigma^{-1} T^k U^\top S^\top S b| \lesssim \gamma^k (\log^c d) \frac{\sqrt{d}}{\sqrt{m}} \|a\|_2 \|\Sigma^{-1}\|_2 \|b\|_2.$$

Since this decays exponentially in $k$ at a rate of $\gamma < 1/2$, the sum of all terms greater than $k$ is bounded by the $k$th term. As long as

$$m \gtrsim \frac{1}{\varepsilon^2} d^{1+\frac{1}{k}} \log^c n, \tag{9}$$

we have $\gamma = \sqrt{\frac{d \log^c n}{m}} < \varepsilon d^{-1/(2k)}/\log^c n$, so that

$$\sum_{k' \geq k} |a^\top V \Sigma^{-1} T^{k'} U^\top S^\top S b| \lesssim \frac{\varepsilon}{\sqrt{d}} \|a\|_2 \|\Sigma^{-1}\|_2 \|b\|_2.$$

On the other hand, by Lemma 9, increasing $m$ by a $C^k$ factor, we have for all $k$ that

$$|a^\top V^\top \Sigma^{-1} U^\top S^\top S (U U^\top S^\top S)^k b| \lesssim \frac{1}{2^k} \frac{\varepsilon}{\sqrt{d}} \|a\|_2 \|b\|_2 \|\Sigma^{-1}\|_2$$

8

with probability at least $1 - 1/\text{poly}(d)$, as long as $m \gtrsim d\log^{2k+2} n \log^2(d/\varepsilon)/\varepsilon^2$. Since the $T^k$ term can be expanded as a sum of $2^k$ terms of this form, we get that

$$\sum_{k'=1}^{k} |a^\top V \Sigma^{-1} T^k U^\top S^\top S b| \lesssim \frac{\varepsilon}{\sqrt{d}} \|a\|_2 \|b\|_2 \|\Sigma^{-1}\|_2$$

with probability at least $1 - 1/\text{poly}(d)$, as long as $m \gtrsim d(C\log n)^{2k+2}\log^2(d/\varepsilon)/\varepsilon^2$ for a sufficiently large constant $C$. Combining with (9), the result holds as long as

$$m \gtrsim \frac{d\log^c n}{\varepsilon^2} \max((C\log n)^{2k+2}, d^{\frac{1}{k}})$$

for any $k$. Setting $k = \Theta(\sqrt{\frac{\log d}{\log\log n}})$ gives the result. ■

**Combining Different Matrices.** In some cases it can make sense to combine different matrices that satisfy the generalization bound.

**Theorem 11.** *Let $A \in \mathbb{R}^{n\times d}$, and let $R \in \mathbb{R}^{m\times r}$ and $S \in \mathbb{R}^{r\times n}$ be drawn from distributions of matrices that are $\varepsilon$-approximate OSEs and satisfy the generalization bound (6). Then $RS$ satisfies the generalization bound with a constant factor loss in failure probability and approximation factor.*

We defer the details to Appendix B.

# 4   Proof of Lemma 9

*Proof.* Each column $S_i$ of the subsampled randomized Hadamard transform has the same distribution as $\sigma_i S_i$, where $\sigma_i$ is a random sign. It also has $\langle S_i, S_i \rangle = 1$ for all $i$ and $|\langle S_i, S_j \rangle| \lesssim \frac{\sqrt{\log(1/\delta)}}{\sqrt{r}}$ with probability $1 - \delta$, for any $\delta$ and $i \neq j$. See, e.g., [LDFU13].

By expanding the following product into a sum, and rearranging terms, we obtain

$$a^\top S^\top S(UU^\top S^\top S)^k b$$
$$= \sum_{i_0, j_0, i_1, j_1, \cdots, i_k, j_k} a_{i_0} b_{j_k} \sigma_{i_0} \sigma_{i_1} \cdots \sigma_{i_k} \sigma_{j_0} \sigma_{j_1} \cdots \sigma_{j_k}$$
$$\cdot \langle S_{i_0}, S_{j_0}\rangle (UU^\top)_{j_0, i_1} \langle S_{i_1}, S_{j_1}\rangle \cdots (UU^\top)_{j_{k-1}, i_k} \langle S_{i_k}, S_{j_k}\rangle$$
$$= \sum_{i_0, j_k} a_{i_0} b_{j_k} \sigma_{i_0} \sigma_{j_k} \sum_{j_0, i_1, j_1, \cdots, i_k} \sigma_{i_1} \cdots \sigma_{i_k} \sigma_{j_0} \sigma_{j_1} \cdots \sigma_{j_{k-1}}$$
$$\cdot \langle S_{i_0}, S_{j_0}\rangle (UU^\top)_{j_0, i_1} \langle S_{i_1}, S_{j_1}\rangle \cdots (UU^\top)_{j_{k-1}, i_k} \langle S_{i_k}, S_{j_k}\rangle$$
$$= \sum_{i_0, j_k} \sigma_{i_0} \sigma_{j_k} Z_{i_0, j_k}$$

where $Z_{i_0, j_k}$ is defined to be

$$Z_{i_0, j_k} = a_{i_0} b_{j_k} \sum_{\substack{i_1, \cdots i_k \\ j_0, \cdots j_{k-1}}} \prod_{c=1}^{k} \sigma_{i_c} \prod_{c=0}^{k-1} \sigma_{j_c} \cdot \prod_{c=0}^{k} \langle S_{i_c}, S_{j_c}\rangle \prod_{c=1}^{k} (UU^\top)_{i_{c-1}, j_c}$$

Note that $Z_{i_0, j_k}$ is independent of $\sigma_{i_0}$ and $\sigma_{j_k}$. We observe that in the above expression if $i_0 = j_0$, $i_1 = j_1$, $\cdots$, $i_k = j_k$, then the sum over these indices equals $a^\top (UU^\top) \cdots (UU^\top) b = 0$, since

9

$\langle S_{i_c}, S_{j_c} \rangle = 1$ in this case for all $c$. Moreover, the sum over all indices conditioned on $i_k = j_k$ is equal to 0. Indeed, in this case, the expression can be factored into the form $\zeta \cdot U^\top b$, for some random variable $\zeta$, but $U^\top b = 0$.

Let $W$ be a matrix with $W_{i,j} = \sigma_i \sigma_j Z_{i,j}$. We need Khintchine's inequality:

**Fact 12** (Khintchine's Inequality). *Let $\sigma_1, \ldots, \sigma_n$ be i.i.d. sign random variables, and let $z_1, \ldots, z_n$ be real numbers. Then there are constants $C, C' > 0$ so that*

$$\Pr \left[ \left| \sum_{i=1}^n z_i \sigma_i \right| \geq Ct\|z\|_2 \right] \leq e^{-C't^2}.$$

We note that Khintchine's inequality sometimes refers to bounds on the moment of $|\sum_i z_i \sigma_i|$, though the above inequality follows readily by applying a Markov bound to the high moments.

We apply Fact 12 to each column of $W$, so that if $W_i$ is the $i$-th column, we have by a union bound that with probability $1 - 1/\mathrm{poly}(n)$, $\|W_i\|_2 = O(\|Z_i\|_2 \sqrt{\log n})$ simultaneously for all columns $i$. It follows that with the same probability, $\|W\|_F^2 = O(\|Z\|_F^2 \log n)$, that is, $\|W\|_F = O(\|Z\|_F \sqrt{\log n})$. We condition on this event in the remainder.

Thus, it remains to bound $\|Z\|_F$. By squaring $Z_{i_0,j_0}$ and using that $\mathbf{E}[\sigma_i \sigma_j] = 1$ if $i = j$ and 0 otherwise, we have,

$$\mathbf{E}_\sigma[Z_{i_0,j_k}^2] = a_{i_0}^2 b_{j_k}^2 \sum_{\substack{i_1,\cdots i_k \\ j_0,\cdots j_{k-1}}} \prod_{c=0}^k \langle S_{i_c}, S_{j_c} \rangle^2 \prod_{c=1}^k (UU^\top)_{i_{c-1},j_c}^2 \tag{10}$$

We defer to Appendix E the proof that

$$\mathbf{E}_S[\|Z\|_F^2] \leq (O(d(\log n)/r) + 1)^{k-1} \cdot (d\|b\|_2^2(\log^2 n)/r^2 + \|b\|_2^2(\log n)/r)$$

Note that we also have the bound:

$$(O(d(\log n)/r) + 1)^{k-1} \leq (e^{O(d(\log n)/r)})^{k-1} \leq e^{O(kd(\log n)/r)} \leq O(1)$$

for any $r = \Omega(kd \log n)$.

Having computed the expectation of $\|Z\|_F^2$, we now would like to show concentration. Consider a specific

$$Z_{i_0,j_k} = a_{i_0} b_{j_k} \sum_{i_k} \sigma_{i_k} \langle S_{i_k}, S_{j_k} \rangle \cdots \sum_{j_1} \sigma_{j_1}(UU^\top)_{j_1,i_2} \sum_{i_1} \sigma_{i_1} \langle S_{i_1}, S_{j_1} \rangle \sum_{j_0} \sigma_{j_0} \langle S_{i_0}, S_{j_0} \rangle (UU^\top)_{j_0,i_1}.$$

By Fact 12, for each fixing of $i_1$, with probability $1 - 1/\mathrm{poly}(n)$, we have

$$\sum_{j_0} \sigma_{j_0} \langle S_{i_0}, S_{j_0} \rangle (UU^\top)_{j_0,i_1} = O(\sqrt{\log n}) \left( \sum_{j_0} \langle S_{i_0}, S_{j_0} \rangle^2 (UU^\top)_{j_0,i_1}^2 \right)^{\frac{1}{2}}. \tag{11}$$

Now, we can apply Khintchine's inequality for each fixing of $j_1$, and combine this with (11). With

probability $1 - 1/\text{poly}(n)$, again we have

$$\sum_{i_1} \sigma_{i_1} \langle S_{i_1}, S_{j_1} \rangle \sum_{j_0} \sigma_{j_0} \langle S_{i_0}, S_{j_0} \rangle (UU^\top)_{j_0, i_1}$$

$$= \sum_{i_1} \sigma_{i_1} \langle S_{i_1}, S_{j_1} \rangle O(\sqrt{\log n}) \left( \sum_{j_0} \langle S_{i_0}, S_{j_0} \rangle^2 (UU^\top)^2_{j_0, i_1} \right)^{\frac{1}{2}}$$

$$= O(\log n) \left( \sum_{i_1} \langle S_{i_1}, S_{j_1} \rangle^2 \sum_{j_0} \langle S_{i_0}, S_{j_0} \rangle^2 (UU^\top)^2_{j_0, i_1} \right)^{\frac{1}{2}}$$

 Thus, we can apply Khintchine's inequality recursively over all the $2k$ indexes $j_0, i_1, j_1, \cdots, j_{k-1}, i_k$, from which it follows that with probability $1 - 1/\text{poly}(n)$, for each such $i_0, j_k$, we have $Z^2_{i_0, j_k} = O(\log^k n) \underset{S}{\mathbf{E}}[Z^2_{i_0, j_k}]$, using (17). We thus have with this probability, that $\|Z\|_F^2 = O(\log^k n) \underset{S}{\mathbf{E}}[\|Z\|_F^2]$, completing the proof. ∎

# 5   Lower bound for $\ell_2$ and $\ell_\infty$ guarantee

We prove a lower bound for the $\ell_2$ guarantee, which immediately implies a lower bound for the $\ell_\infty$ guarantee.

**Definition 13.** *Given a matrix $A \in \mathbb{R}^{n \times d}$, vector $b \in \mathbb{R}^n$ and matrix $S \in \mathbb{R}^{r \times n}$, denote $x^* = A^\dagger b$. We say that an algorithm $\mathcal{A}(A, b, S)$ that outputs a vector $x' = (SA)^\dagger Sb$ "succeeds" if the following property holds: $\|x' - x^*\|_2 \lesssim \varepsilon \|b\|_2 \cdot \|A^\dagger\|_2 \cdot \|Ax^* - b\|_2$.*

**Theorem 14.** *Suppose $\Pi$ is a distribution over $\mathbb{R}^{m \times n}$ with the property that for any $A \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$, $\underset{S \sim \Pi}{\Pr}[\mathcal{A}(A, b, S) \text{ succeeds }] \geq 19/20$. Then $m \gtrsim \min(n, d/\varepsilon^2)$.*

*Proof.* The proof uses Yao's minimax principle. Let $\mathcal{D}$ be an arbitrary distribution over $\mathbb{R}^{n \times (d+1)}$, then $\underset{(A,b) \sim \mathcal{D}}{\mathbf{E}} \underset{S \sim \Pi}{\mathbf{E}}[\mathcal{A}(A, b, S) \text{ succeeds }] \geq 1 - \delta$. Switching the order of probabilistic quantifiers, an averaging argument implies the existence of a fixed matrix $S_0 \in \mathbb{R}^{m \times n}$ such that

$$\underset{(A,b) \sim \mathcal{D}}{\mathbf{E}}[\mathcal{A}(A, b, S_0) \text{ succeeds }] \geq 1 - \delta.$$

Thus, we must construct a distribution $\mathcal{D}_{\text{hard}}$ such that

$$\underset{(A,b) \sim \mathcal{D}_{\text{hard}}}{\mathbf{E}}[\mathcal{A}(A, b, S_0) \text{ succeeds }] \geq 1 - \delta,$$

cannot hold for any $\Pi_0 \in \mathbb{R}^{m \times n}$ which does not satisfy $m = \Omega(d/\varepsilon^2)$. The proof can be split into three parts. First, we prove a useful property. Second, we prove a lower bound for the case $\text{rank}(S) \geq d$. Third, we show why $\text{rank}(S) \geq d$ is necessary.

   (I) We show that $[SA, Sb]$ are independent Gaussian, if both $[A, b]$ and $S$ are orthonormal matrices. We can rewrite $SA$ in the following sense,

$$\underbrace{S}_{m \times n} \cdot \underbrace{A}_{n \times d} = \underbrace{S}_{m \times n} \underbrace{R}_{n \times n} \underbrace{R^\top}_{n \times n} \underbrace{A}_{n \times d} = S \begin{bmatrix} S^\top & \overline{S}^\top \end{bmatrix} \begin{bmatrix} S \\ \overline{S} \end{bmatrix} A = \begin{bmatrix} I_m & 0 \end{bmatrix} \begin{bmatrix} S \\ \overline{S} \end{bmatrix} A = \begin{bmatrix} I_m & 0 \end{bmatrix} \underbrace{\widetilde{A}}_{n \times d} = \underbrace{\widetilde{A}_m}_{m \times d}$$

11

where $\overline{S}$ is the complement of the orthonormal basis $S$, $I_m$ is a $m \times m$ identity matrix, and $\widetilde{A}_m$ is the left $m \times d$ submatrix of $\widetilde{A}$. Thus, using [J$^+$06] as long as $m = o(\sqrt{n})$ (because of $n = \Omega(d^3)$) the total variation distance between $[SA, Sb]$ and a random Gaussian matrix is small, i.e.,

$$D_{TV}([SA, Sb], H) \leq 0.01 \tag{12}$$

where each entry of $H$ is i.i.d. Gaussian $\mathcal{N}(0, 1/n)$.

(II) Here we prove the theorem in the case when $S$ has rank $r \geq d$ (we will prove this is necessary in part III. Writing $S = U\Sigma V^\top$ in its SVD, we have

$$\underbrace{S}_{m \times n} A = \underbrace{U}_{m \times r} \underbrace{\Sigma}_{r \times r} \underbrace{V^\top}_{r \times n} RR^\top A = U\Sigma G \tag{13}$$

where $R = \begin{bmatrix} V & \overline{V} \end{bmatrix}$. By a similar argument in Equation (12), as long as $r = o(\sqrt{n})$ we have that $G$ also can be approximated by a Gaussian matrix, where each entry is sampled from i.i.d. $\mathcal{N}(0, 1/n)$. Similarly, $Sb = U\Sigma h$, where $h$ also can be approximated by a Gaussian matrix, where each entry is sampled from i.i.d. $\mathcal{N}(0, 1/n)$.

Since $U$ has linearly independent columns, $(U\Sigma G)^\dagger U\Sigma h = (\Sigma G)^\dagger U^\top U\Sigma h = (\Sigma G)^\dagger \Sigma h$.

The $r \times d$ matrix $G$ has $SVD$ $G = \underbrace{R}_{r \times d} \underbrace{\widetilde{\Sigma}}_{d \times d} \underbrace{T}_{d \times d}$, and applying the pseudo-inverse property again, we have

$$\|(SA)^\dagger Sb\|_2 = \|(\Sigma G)^\dagger \Sigma h\|_2 = \|(\Sigma R\widetilde{\Sigma} T)^\dagger \Sigma h\|_2 = \|T^\dagger (\Sigma R\widetilde{\Sigma})^\dagger \Sigma h\|_2 = \|(\Sigma R\widetilde{\Sigma})^\dagger \Sigma h\|_2$$
$$= \|\widetilde{\Sigma}^\dagger (\Sigma R)^\dagger \Sigma h\|_2,$$

where the the first equality follows by Equation (13), the second equality follows by the $SVD$ of $G$, the third and fifth equality follow by properties of the pseudo-inverse[1] when $T$ has orthonormal rows and $\widetilde{\Sigma}$ is a diagonal matrix, and the fourth equality follows since $\|T^\dagger\|_2 = 1$ and $T$ is an orthonormal basis.

Because each entry of $G = R\widetilde{\Sigma} T \in \mathbb{R}^{r \times d}$ is sampled from an i.i.d. Gaussian $\mathcal{N}(0, 1)$, using the result of [Ver10] we can give an upper bound for the maximum singular value of $G$: $\|\widetilde{\Sigma}\| \lesssim \sqrt{\frac{r}{n}}$ with probability at least .99. Thus,

$$\|\widetilde{\Sigma}^\dagger (\Sigma R)^\dagger \Sigma h\|_2 \geq \sigma_{\min}(\widetilde{\Sigma}^\dagger) \cdot \|(\Sigma R)^\dagger \Sigma h\|_2 = \frac{1}{\sigma_{\max}(\widetilde{\Sigma})} \|(\Sigma R)^\dagger \Sigma h\|_2 \gtrsim \sqrt{n/r} \|(\Sigma R)^\dagger \Sigma h\|_2.$$

Because $h$ is a random Gaussian vector which is independent of $(\Sigma R)^\dagger \Sigma$, by Claim 15, $\mathbf{E}_h[\|(\Sigma R)^\dagger \Sigma h\|_2^2] = \frac{1}{n} \cdot \|(\Sigma R)^\dagger \Sigma\|_F^2$, where each entry of $h$ is sampled from i.i.d. Gaussian $\mathcal{N}(0, 1/n)$. Then, using the Pythagorean Theorem,

$$\|(\Sigma R)^\dagger \Sigma\|_F^2 = \|(\Sigma R)^\dagger \Sigma RR^\top\|_F^2 + \|(\Sigma R)^\dagger \Sigma (I - RR^\top)\|_F^2$$
$$\geq \|(\Sigma R)^\dagger \Sigma RR^\top\|_F^2$$
$$= \|(\Sigma R)^\dagger \Sigma R\|_F^2$$
$$= \text{rank}(\Sigma R)$$
$$= \text{rank}(SA)$$
$$= d.$$

---

[1] https://en.wikipedia.org/wiki/Moore-Penrose_pseudoinverse

Thus, $\|x' - x^*\|_2 \gtrsim \sqrt{d/r} \geq \sqrt{d/m} = \varepsilon$.

(III) Now we show that we can assume that $\text{rank}(S) \geq d$.

We sample $A, b$ based on the following distribution $\mathcal{D}_{\text{hard}}$: with probability $1/2$, $A, b$ are sampled from $\mathcal{D}_1$; with probability $1/2$, $A, b$ are sampled from $\mathcal{D}_2$. In distribution $\mathcal{D}_1$, $A$ is a random orthonormal basis and $d$ is always orthogonal to $A$. In distribution $\mathcal{D}_2$, $A$ is a $d \times d$ identity matrix in the top-$d$ rows and 0s elsewhere, while $b$ is a random unit vector. Then, for any $(A, b)$ sampled from $\mathcal{D}_1$, $S$ needs to work with probability at least $9/10$. Also for any $(A, b)$ sampled from $\mathcal{D}_2$, $S$ needs to work with probability at least $9/10$. The latter two statements follow since overall $S$ succeeds on $\mathcal{D}_{\text{hard}}$ with probability at least $19/20$.

Consider the case where $A, b$ are sampled from distribution $\mathcal{D}_2$. Then $x^* = b$ and $\text{OPT} = 0$. Then consider $x'$ which is the optimal solution to $\min_x \|SAx - Sb\|_2^2$, so $x' = (SA)^\dagger Sb = (S_L)^\dagger S_L b$, where $S$ can be decomposed into two matrices $S_L \in \mathbb{R}^{r \times d}$ and $S_R \in \mathbb{R}^{r \times (n-d)}$, $S = \begin{bmatrix} S_L & S_R \end{bmatrix}$. Plugging $x'$ into the original regression problem, $\|Ax' - b\|_2^2 = \|A(S_L)^\dagger S_L b - b\|_2^2$, which is at most $(1 + \varepsilon) \text{OPT} = 0$. Thus $\text{rank}(S_L)$ is $d$. Since $S_L$ is a submatrix of $S$, the rank of $S$ is also $d$. ∎

It remains to define several tools which are used in the main proof of the lower bound.

**Claim 15.** *For any matrix $A \in \mathbb{R}^{n \times d}$, if each entry of a vector $g \in \mathbb{R}^d$ is chosen from an i.i.d Gaussian $\mathcal{N}(0, \sigma^2)$, then $\mathbf{E}_g[\|Ag\|_2^2] = \sigma^2 \|A\|_F^2$ .*

*Proof.*

$$
\begin{aligned}
\mathbf{E}_g[\|Ag\|_2^2] &= \mathbf{E}_g\left[\sum_{i=1}^n (\sum_{j=1}^d A_{ij} g_j)^2\right] \\
&= \mathbf{E}_g\left[\sum_{i=1}^n (\sum_{j=1}^d A_{ij}^2 g_j^2 + \sum_{j \neq j'} A_{ij} A_{ij'} g_j g_{j'})\right] \\
&= \sum_{i=1}^n \sum_{j=1}^d A_{ij}^2 \sigma^2 \\
&= \sigma^2 \|A\|_F^2.
\end{aligned}
$$

∎

Let $g_1, g_2, \cdots, g_t$ be i.i.d. $\mathcal{N}(0, 1)$ random variables. The random variables $\sum_{i=1}^t g_i^2$ are $\mathcal{X}^2$ with $t$ degree of freedom. Furthermore, the following tail bounds are known.

**Fact 16** (Lemma 1 of [LM00]). *Let $g_1, g_2, \cdots, g_t$ be i.i.d. $\mathcal{N}(0, 1)$ random variables. Then for any $x \geq 0$,*

$$
\Pr\left[\sum_{i=1}^t g_i^2 \geq t + 2\sqrt{tx} + 2x\right] \leq \exp(-x),
$$

*and*

$$
\Pr\left[\sum_{i=1}^t g_i^2 \leq t - 2\sqrt{tx}\right] \leq \exp(-x).
$$

**Definition 17.** *Given a matrix $A \in \mathbb{R}^{n \times d}$, vector $b \in \mathbb{R}^n$ and matrix $S \in \mathbb{R}^{r \times n}$, denote $x^* = A^\dagger b$. We say that an algorithm $\mathcal{B}(A, b, S)$ that outputs a vector $x' = (SA)^\dagger Sb$ "succeeds" if the following property holds:*

$$\|x' - x^*\|_\infty \lesssim \frac{\varepsilon}{\sqrt{d}} \|b\|_2 \cdot \|A^\dagger\|_2 \cdot \|Ax^* - b\|_2.$$

Applying $\|x' - x\|_\infty \geq \frac{1}{\sqrt{d}} \|x' - x\|_2$ to Theorem 14 ,we obtain the $\ell_\infty$ lower bound as a corollary,

**Corollary 18.** *Suppose $\Pi$ is a distribution over $\mathbb{R}^{m \times n}$ with the property that for any $A \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$,*

$$\Pr_{S \sim \Pi}[\mathcal{B}(A, b, S) \text{ succeeds }] \geq 9/10.$$

*Then $m \gtrsim \min(n, d/\varepsilon^2)$.*

# References

[ADL+12]   Patrick Amestoy, Iain S. Duff, Jean-Yves L'Excellent, Yves Robert, François-Henry Rouet, and Bora Uçar. On computing inverse entries of a sparse matrix in an out-of-core environment. *SIAM J. Scientific Computing*, 34(4), 2012.

[AMT10]   Haim Avron, Petar Maymounkov, and Sivan Toledo. Blendenpik: Supercharging lapack's least-squares solver. *SIAM J. Scientific Computing*, 32(3):1217–1236, 2010.

[BDN15]   Jean Bourgain, Sjoerd Dirksen, and Jelani Nelson. Toward a unified theory of sparse dimensionality reduction in euclidean space. In *STOC*, 2015.

[CMM15]   Michael B Cohen, Cameron Musco, and Christopher Musco. Ridge leverage scores for low-rank approximation. *arXiv preprint arXiv:1511.07263*, 2015.

[CNW16]   Michael B Cohen, Jelani Nelson, and David P. Woodruff. Optimal approximate matrix product in terms of stable rank. In *Proceedings of the 43rd International Colloquium on Automata, Languages and Programming (ICALP 2016), Rome, Italy, July 12-15, arXiv preprint arXiv:1507.02268*, 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.

[CW90]   Don Coppersmith and Shmuel Winograd. Matrix multiplication via arithmetic progressions. *J. Symb. Comput.*, 9(3):251–280, 1990.

[CW09]   Kenneth L. Clarkson and David P. Woodruff. Numerical linear algebra in the streaming model. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*, pages 205–214. ACM, 2009.

[CW13]   Kenneth L Clarkson and David P. Woodruff. Low rank approximation and regression in input sparsity time. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 81–90. ACM, 2013.

[DMIMW12]   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.

[DMMS11]   Petros Drineas, Michael W Mahoney, S Muthukrishnan, and Tamás Sarlós. Faster least squares approximation. *Numerische Mathematik*, 117(2):219–249, 2011.

[Gal14]   François Le Gall. Powers of tensors and fast matrix multiplication. In *International Symposium on Symbolic and Algebraic Computation, ISSAC '14, Kobe, Japan, July 23-25, 2014*, pages 296–303, 2014.

[J+06]   Tiefeng Jiang et al. How many entries of a typical orthogonal matrix can be approximated by independent normals? *The Annals of Probability*, 34(4):1497–1529, 2006.

[LAKD08]   Song Li, Shaikh S. Ahmed, Gerhard Klimeck, and Eric Darve. Computing entries of the inverse of a sparse matrix using the FIND algorithm. *J. Comput. Physics*, 227(22):9408–9427, 2008.

[LDFU13]   Yichao Lu, Paramveer Dhillon, Dean Foster, and Lyle Ungar. Faster ridge regression via the subsampled randomized hadamard transform. In *Proceedings of the Neural Information Processing Systems (NIPS) Conference*, 2013.

[Lee12]    Jeff Leek. Prediction: the lasso vs. just using the top 10 predictors. [http://simplystatistics.org/2012/02/23/prediction-the-lasso-vs-just-using-the-top-10/](http://simplystatistics.org/2012/02/23/prediction-the-lasso-vs-just-using-the-top-10/), 2012.

[LM00]     Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000.

[LMP13]    Mu Li, Gary L Miller, and Richard Peng. Iterative row sampling. In *Foundations of Computer Science (FOCS), 2013 IEEE 54th Annual Symposium on*, pages 127–136. IEEE, 2013.

[MM13]     Xiangrui Meng and Michael 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, 2013.

[NN13]     Jelani Nelson and Huy L Nguyên. Osnap: Faster numerical linear algebra algorithms via sparser subspace embeddings. In *Foundations of Computer Science (FOCS), 2013 IEEE 54th Annual Symposium on*, pages 117–126. IEEE, 2013.

[NN14]     Jelani Nelson and Huy L. Nguyên. Lower bounds for oblivious subspace embeddings. In *Automata, Languages, and Programming - 41st International Colloquium, ICALP 2014, Copenhagen, Denmark, July 8-11, 2014, Proceedings, Part I*, pages 883–894, 2014.

[PKB14]    Dimitris Papailiopoulos, Anastasios Kyrillidis, and Christos Boutsidis. Provable deterministic leverage score sampling. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 997–1006. ACM, 2014.

[RYP03]    Ryan Rifkin, Gene Yeo, and Tomaso Poggio. Regularized least-squares classification. *Nato Science Series Sub Series III Computer and Systems Sciences*, 190:131–154, 2003.

[Sar06]    Tamás Sarlós. Improved approximation algorithms for large matrices via random projections. In *Foundations of Computer Science, 2006. FOCS'06. 47th Annual IEEE Symposium on*, pages 143–152. IEEE, 2006.

[Ver10]    Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.

[Wil12]    Virginia Vassilevska Williams. Multiplying matrices faster than coppersmith-winograd. In *Proceedings of the 44th Symposium on Theory of Computing Conference, STOC 2012, New York, NY, USA, May 19 - 22, 2012*, pages 887–898, 2012.

[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.

[ZF13]     Anastasios Zouzias and Nikolaos M. Freris. Randomized extended kaczmarz for solving least squares. *SIAM J. Matrix Analysis Applications*, 34(2):773–793, 2013.

[ZP04]     Peng Zhang and Jing Peng. Svm vs regularized least squares classification. In *ICPR (1)*, pages 176–179, 2004.

# Appendix

## A   Proof for Gaussian case

**Lemma 19.** *If the entries of $S \in \mathbb{R}^{m \times n}$ are i.i.d. $N(0, 1/m)$, $m = O(d/\varepsilon^2)$, and $U^\top b = 0$, then*

$$|a^\top (SA)^\dagger Sb| \lesssim \frac{\varepsilon \sqrt{\log d}}{\sqrt{d}} \|a\|_2 \|b\|_2 \|\Sigma^{-1}\|_2$$

*for any vectors $a, b$ with probability $1 - 1/\mathrm{poly}(d)$.*

*Proof.* With probability 1, the matrix $SA$ has linearly independent columns, and so $(SA)^\dagger$ is

$$
\begin{aligned}
&= (A^\top S^\top SA)^{-1} A^\top S^\top \\
&= (V\Sigma U^\top S^\top SU\Sigma V^\top)^{-1} V\Sigma U^\top S^\top \\
&= V\Sigma^{-1}(U^\top S^\top SU)^{-1}\Sigma^{-1} V^\top V\Sigma U^\top S^\top \\
&= V\Sigma^{-1}(U^\top S^\top SU)^{-1} U^\top S^\top.
\end{aligned}
$$

Hence, we would like to bound

$$X = a^\top V\Sigma^{-1}(U^\top S^\top SU)^{-1} U^\top S^\top Sb.$$

It is well-known (stated, for example, explicitly in Theorem 2.3 of [Woo14]) that with probability $1 - \exp(-d)$, the singular values of $SU$ are $(1 \pm \varepsilon)$ for $m = O(d/\varepsilon^2)$. We condition on this event. It follows that

$$
\begin{aligned}
& \|V\Sigma^{-1}(U^\top S^\top SU)^{-1} U^\top S\|_2 \\
=\ & \|\Sigma^{-1}(U^\top S^\top SU)^{-1} U^\top S\|_2 \\
\leq\ & \|\Sigma^{-1}\|_2 \|(U^\top S^\top SU)^{-1}\|_2 \|U^\top S\|_2 \\
\leq\ & \|\Sigma^{-1}\|_2 \cdot \frac{1}{1-\varepsilon} \cdot (1+\varepsilon) \\
=\ & O(\|\Sigma^{-1}\|_2),
\end{aligned}
$$

where the first equality uses that $V$ is a rotation, the first inequality follows by sub-multiplicativity, and the second inequality uses that the singular values of $SU$ are in the range $[1-\varepsilon, 1+\varepsilon]$. Hence, with probability $1 - \exp(-d)$,

$$\|a^\top V\Sigma^{-1}(U^\top S^\top SU)^{-1} U^\top S^\top\|_2 = O(\|\Sigma^{-1}\|_2 \|a\|_2). \tag{14}$$

The main observation is that since $U^\top b = 0$, $SU$ is statistically independent from $Sb$. Hence, $Sb$ is distributed as $N(0, \|b\|_2^2 I_m)$, conditioned on the vector $a^\top V\Sigma^{-1}(U^\top S^\top SU)^{-1} U^\top S^\top$. It follows that conditioned on the value of $a^\top V\Sigma^{-1}(U^\top S^\top SU)^{-1} U^\top S^\top$, $X$ is distributed as

$$N(0, \|b\|_2^2 \|a^\top V\Sigma^{-1}(U^\top S^\top SU)^{-1} U^\top S^\top\|_2^2 / m),$$

and so using (14) , with probability $1 - 1/\mathrm{poly}(d)$, we have $|X| = O(\varepsilon\sqrt{\log d}\|a\|_2 \|b\|_2 \|\Sigma^{-1}\|_2 / \sqrt{d})$. $\blacksquare$

# B Combining Different Matrices

In some cases it can make sense to combine different matrices that satisfy the generalization bound.

**Theorem 11.** *Let $A \in \mathbb{R}^{n \times d}$, and let $R \in \mathbb{R}^{m \times r}$ and $S \in \mathbb{R}^{r \times n}$ be drawn from distributions of matrices that are $\varepsilon$-approximate OSEs and satisfy the generalization bound* (6). *Then $RS$ satisfies the generalization bound with a constant factor loss in failure probability and approximation factor.*

*Proof.* For any vectors $a, b$, and $x^* = A^\dagger b$ we want to show

$$|a^\top (RSA)^\dagger RSb - a^\top x^*| \lesssim \frac{\varepsilon}{\sqrt{d}} \|a\|_2 \|b - Ax^*\|_2 \|A^\dagger\|_2$$

As before, it suffices to consider the $x^* = 0$ case. We have with probability $1 - \delta$ that

$$|a^\top (SA)^\dagger Sb| \lesssim \frac{\varepsilon}{\sqrt{d}} \|a\|_2 \|b\|_2 \|A^\dagger\|_2;$$

suppose this happens. We also have by the properties of $R$, applied to $SA$ and $Sb$, that

$$|a^\top (RSA)^\dagger RSb - a^\top (SA)^\dagger Sb| \lesssim \frac{\varepsilon}{\sqrt{d}} \|a\|_2 \|Sb\|_2 \|(SA)^\dagger\|_2.$$

Because $S$ is an OSE, we have $\|Sb\|_2 \leq (1 + \varepsilon)$ and $\|(SA)^\dagger\|_2 \gtrsim (1 - \varepsilon)\|A^\dagger\|_2$. Therefore

$$|a^\top (RSA)^\dagger RSb| \lesssim \frac{\varepsilon}{\sqrt{d}} \|a\|_2 \|b\|_2 \|A^\dagger\|_2$$

$\blacksquare$

We describe a few of the applications of combining sketches.

## B.1 Removing dependence on $n$ via Count-Sketch

One of the limitations of the previous section is that the choice of $k$ depends on $n$. To prove that theorem, we have to assume that $\log d > \log \log n$. Here, we show an approach to remove that assumption.

The main idea is instead of applying matrix $S \in \mathbb{R}^{m \times n}$ to matrix $A \in \mathbb{R}^{n \times d}$ directly, we pick two matrices $S \in \mathbb{R}^{m \times \text{poly}(d)}$ and $C \in \mathbb{R}^{\text{poly}(d) \times n}$, e.g. $S$ is FastJL matrix and $C$ is Count-Sketch matrix with $s = 1$. We first compute $C \cdot A$, then compute $S \cdot (CA)$. The benefit of these operations is $S$ only needs to multiply with a matrix $(CA)$ that has $\text{poly}(d)$ rows, thus the assumption we need is $\log d > \log \log(\text{poly}(d))$ which is always true. The reason for choosing $C$ as a Count-Sketch matrix with $s = 1$ is: (1) $\text{nnz}(CA) \leq \text{nnz}(A)$ (2) The running time is $O(\text{poly}(d) \cdot d + \text{nnz}(A))$.

## B.2 Combining Gaussians and SRHT

By combining Gaussians with SRHT matrices, we can embed into the optimal dimension $O(d/\varepsilon^2)$ with fast $\widetilde{O}(nd \log n + d^\omega / \varepsilon^4)$ embedding time.

## B.3 Combining all three

By taking Gaussians times SRHT times Count-Sketch, we can embed into the optimal dimension $O(d/\varepsilon^2)$ with fast $O(\text{nnz}(A) + d^4 \text{poly}(\frac{1}{\varepsilon}, \log d))$ embedding time.

# C  Count-Sketch does not obey the $\ell_\infty$ guarantee

Here we demonstrate an $A$ and a $b$ such that Count-Sketch will not satisfy the $\ell_\infty$ guarantee with constant probability, so such matrices cannot satisfy the generalization guarantee (6) with high probability.

**Theorem 20.** *Let $S \in \mathbb{R}^{m \times n}$ be drawn as a Count-Sketch matrix with $s$ nonzeros per column. There exists a matrix $A \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$ such that, if $s^2 d \lesssim m \lesssim \sqrt{d^3 s}$, then the "true" solution $x^* = A^\dagger b$ and the approximation $x' = (SA)^\dagger Sb$ have large $\ell_\infty$ distance with constant probability:*

$$\|x' - x^*\|_\infty \gtrsim \sqrt{\frac{d}{ms}} \|b\|_2.$$

*Plugging in $m = d^{1.5}$ and $s = d^{0.25}$ we find that*

$$\|x' - x^*\|_\infty \gtrsim 1/d^{3/8} \|b\|_2 \gg 1/\sqrt{d} \|b\|_2,$$

*even though such a matrix is an OSE with probability exponential in $s$. Therefore there exists a constant $c$ for which this matrix does not satisfy the generalization guarantee (6) with $1 - \frac{c}{d}$ probability.*

*Proof.* We choose the matrix $A$ to be the identity on its top $d$ rows: $A = \begin{bmatrix} I_d \\ 0 \end{bmatrix}$. Choose some $\alpha \geq 1$, set the value of the first $d$ coordinates of vector $b$ to be $\frac{1}{\sqrt{d}}$ and set the value to be $1/\sqrt{\alpha}$ for the next $\alpha$ coordinates, with the remaining entries all zero. Note that $\|b\|_2 = \sqrt{2}$, $x^* = (1/\sqrt{d}, \ldots, 1/\sqrt{d})$, and $\|Ax^* - b\|_2 = 1$.

Let $S_k$ denote the $k$th column vector of matrix $S \in \mathbb{R}^{m \times n}$. We define two events, Event I, $\forall k' \in [d]$ and $k' \neq k$, we have $\text{supp}(S_{k'}) \cap \text{supp}(S_k) = \emptyset$; Event II, $\exists$ a unique $k' \in \{d+1, d+2, \cdots, d+\alpha\}$ such that $|\text{supp}(S_{k'}) \cap \text{supp}(S_k)| = 1$, and all other $k'$ have $\text{supp}(S_{k'}) \cap \text{supp}(S_k) = \emptyset$. Using Claim 21, with probability at least .99 there exists a $k$ for which both events hold.

Given the constructions of $A$ and $b$ described early, it is obvious that

$$Ax - b = \left[ x_1 - \tfrac{1}{\sqrt{d}}, \cdots, x_d - \tfrac{1}{\sqrt{d}}, -\tfrac{1}{\sqrt{\alpha}}, \cdots, -\tfrac{1}{\sqrt{\alpha}}, 0, \cdots, 0 \right]^\top.$$

Conditioned on event I and II are holding, then denote $\text{supp}(S_j) = \{i_1, i_2, \cdots, i_s\}$. Consider the terms involving $x_j$ in the quadratic form

$$\min_x \|SAx - Sb\|_2^2.$$

it can be written as $(s-1)(x_j - 1/\sqrt{d})^2 + (x_j - 1/\sqrt{d} \pm 1/\sqrt{\alpha})^2$. Hence the optimal $x'$ will have $x'_j = \frac{1}{\sqrt{d}} \pm \frac{1}{s\sqrt{\alpha}}$, which is different from the desired $1/\sqrt{d}$ by $\frac{1}{s\sqrt{\alpha}}$. Plugging in our requirement of $\alpha \approx m^2/(s^3 d^2)$, we have

$$\|x' - x^*\|_\infty \geq \frac{1}{s\sqrt{\alpha}} \gtrsim c\sqrt{\frac{sd^2}{m^2}} \gtrsim \frac{1}{\sqrt{d}}$$

where the last inequality follows by $m \lesssim \sqrt{sd^3}$. Thus, we get the result. ∎

**Claim 21.** *If $m = \Omega(s^2 d)$, $m = o(d^2)$, $\alpha < d$, and $\alpha = O(\frac{m^2}{s^3 d^2})$, with probability at least .99 there exists a $k \in [d]$ for which both event I and II hold.*

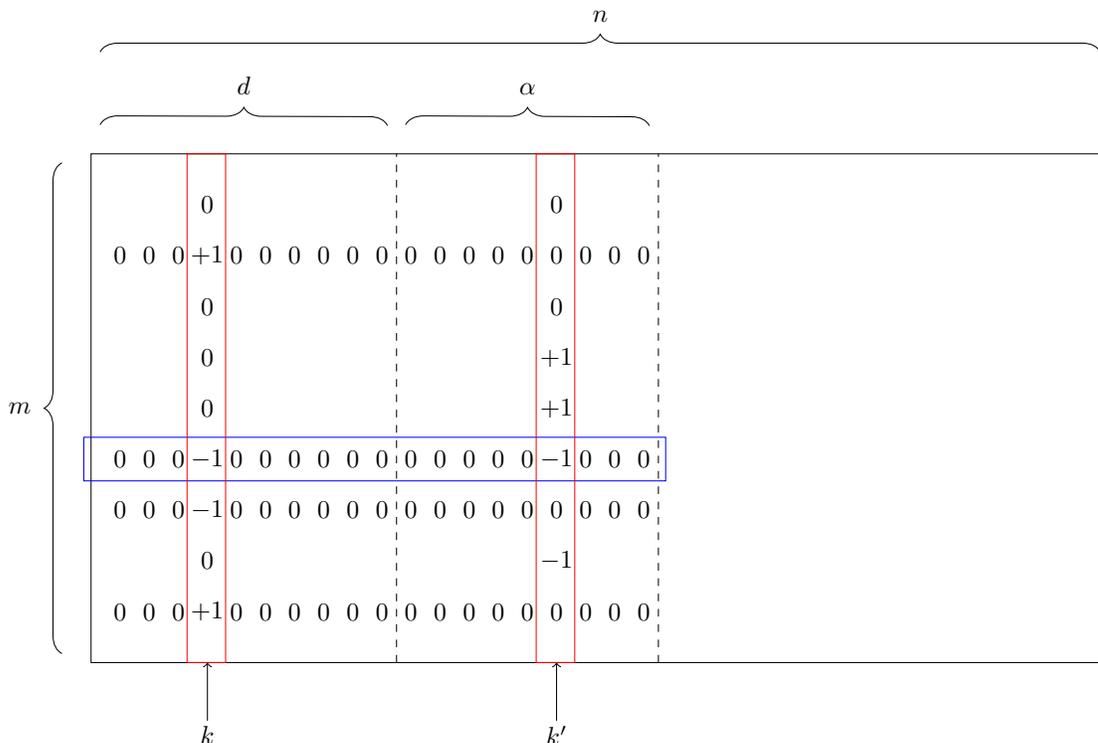Figure 2: Count Sketch matrix $S \in \mathbb{R}^{m \times n}$. Event I, for any $k' \in [d]$ and $k' \neq k$, $\text{supp}(S_k) \cap \text{supp}(S_{k'}) = \emptyset$. Event II, there exists a unique $k' \in \{d+1, d+2, \cdots, d+\alpha\}$ such that $S_k$ and $S_{k'}$ intersect at exactly one location(row index).

*Proof.* If $m = \Omega(s^2 d)$, then for any $i$ in $\{1, 2, ..., d\}$, let $X_i$ be an indicator that the entries of column i are disjoint from all $i'$ in $[d] \backslash \{i\}$. Then $\mathbf{E}[X_i] \geq .9999$, so by Markov's inequality, with probability .99, we have $.99d$ columns having this property (indeed, the expected value of $d - X$ is at most $.0001d$, so $\Pr[d - X \geq .01d] \leq \frac{\mathbf{E}[d-X]}{.01d} \leq \frac{.0001d}{.01d} = .01$). Define Event $E$ to be that $.99d$ columns of first $d$ columns have the property that the entries of that column are disjoint from all the other $d - 1$ columns. Let $S$ be the set of these $.99d$ columns. Let $N$ be the union of supports of columns in $S$.

Each column $i$ in $\{d+1, ..., d+\alpha\}$ chooses $s$ non-zero entries. Define event $F($ which is similar as event $E)$ to be that $.99\alpha$ columns of the next $\alpha$ columns have the property that the entries of that column are disjoint from all the other $\alpha - 1$ columns. By the same argument, since $\alpha < d$, with probability .99, we have $.99\alpha$ columns in $\{d+1, ..., d+\alpha\}$ being disjoint from other columns in $\{d+1, ..., d+\alpha\}$. Condition on event $F$ holding. Let $L$ be the multiset union of supports of all columns in $\{d+1, ..., d+\alpha\}$. Then $L$ has size $\alpha \cdot s$. Let $M$ be the union of supports of all columns in $\{d+1, ..., d+\alpha\}$, that is, the set union rather than the multiset union. Note that $|M| \geq .99\alpha \cdot s$ because of $.99\alpha$ columns are disjoint from each other.

The intersection size $x$ of $N$ and $M$ is hyper-geometrically distributed with expectation

$$\mathbf{E}[x] = \frac{s|S| \cdot |M|}{m}.$$

By a lower tail bound for the hypergeometric distribution [2] ,

$$\Pr[x \leq (p-t)n] \leq \exp(-2t^2 n),$$

where $p = s \cdot |S|/m$ and $n = |M|$, so

$$\Pr[x \leq \mathbf{E}[x] - t \cdot |M|] \leq \exp(-2t^2 \cdot |M|) \leq 0.01,$$

where the last inequality follows by setting $t = \Theta(1/\sqrt{|M|})$. Thus, we get with probability .99, the intersection size is at least $\frac{s|S| \cdot |M|}{m} - \Theta(\sqrt{|M|})$ .

Now let $W$ be the distinct elements in $L \backslash M$, so necessarily $|W| \leq .01\alpha \cdot s$. By an upper tail bound for the hypergeometric distribution, the intersection size $y$ of $W$ and $N$ satisfies

$$\Pr[y \geq (p+t)n] \leq \exp(-2t^2 n),$$

where $p = s \cdot |S|/m$ and $n = |W|$, we again get

$$\Pr[y \geq \mathbf{E}[y] + t \cdot |W|] \leq \exp(-2t^2 \cdot |W|).$$

If $|W| = 0$, then $y = 0$. Otherwise, we can set $t = \Theta(1/\sqrt{|W|})$ so that this probability is less than .01, and we get with probability .99, the intersection size $y$ is at most $s \cdot |S| \cdot |W|/m + \Theta(\sqrt{|W|})$. Note that we have that $\Theta(\sqrt{|M|})$ and $\Theta(\sqrt{|W|})$ are bounded by $\Theta(\sqrt{s \cdot \alpha})$. Setting $\alpha = O(\frac{m^2}{s^3 d^2})$ suffices to ensure $y$ is at most $(1.01)s \cdot |S| \cdot |W|/m$, and earlier that $x$ is at least $.99 \cdot s \cdot |S| \cdot |M|/m$.

The probability one of the $|S|$ blocks in $N$ has two or more intersections with $M$ is less than $\binom{x}{2}$ times the probability two random distinct items in the intersection land in the block. This probability is

$$\frac{\binom{x}{2} \cdot \binom{s}{2}}{\binom{s \cdot |S|}{2}} = \Theta(x^2/|S|^2) = \Theta(x^2/d^2) = \Theta(m^2/(d^4 s^2)).$$

So the expected number of such blocks is $\Theta(m^2 sd/(d^4 s^2)) = \Theta(m^2/(d^3 s))$ which is less than $(.99 \cdot s \cdot |S| \cdot |M|)/(2m) \leq X/2$ if $m = o(d^2)$, which we have. So, there are at least $x/2$ blocks which have intersection size exactly 1 with $N$. Note that the number of intersections of the $|S|$ blocks with $W$ is at most $y$, which is at most $(1.01)s \cdot |S| \cdot |W|/m \leq (1.01)s \cdot |S| \cdot \frac{1}{99} \cdot |M|/m < x/2$, and therefore there exists a block, that is, a column among the first $d$ columns, which intersects $M$ in exactly one position and does not intersect $W$. This is our desired column. Thus, we complete the proof. ■

# D    Leverage score sampling does not obey the $\ell_\infty$ guarantee

Not only does Count-Sketch fail, but so does leverage score sampling, which is a technique that takes a subsample of rows of $A$ with rescaling. In this section we show an $A$ and a $b$ such that leverage score sampling will not satisfy the $\ell_\infty$ guarantee. We start with a formal definition of leverage scores.

**Definition 22** (Leverage Scores). *Given an arbitrary $n \times d$ matrix $A$, with $n > d$, let $U$ denote the $n \times d$ matrix consisting of the $d$ left singular vectors of $A$, let $U_{(i)}$ denote the $i-$th row of the matrix $U$, so $U_{(i)}$ is a row vector. Then the leverage scores of the rows of $A$ are given by $l_i = \|U_{(i)}\|_2^2$, for $i \in [n]$.*

---

[2]https://en.wikipedia.org/wiki/Hypergeometric_distribution

The leverage score sampling matrix can be thought of as a square diagonal matrix $D \in \mathbb{R}^{n \times n}$ with diagonal entries chosen from some distribution. If $D_{ii} = 0$, it means we do not choose the $i$-th row of matrix $A$., If $D_{ii} > 0$, it means we choose that row of the matrix $A$ and also rescale that row. We show that the leverage score sampling matrix cannot achieve $\ell_\infty$ guarantee, nor can it achieve our notion of generalization error.

**Theorem 23.** *Let $D \in \mathbb{R}^{n \times n}$ be a leverage score sampling matrix with $m$ nonzeros on the diagonal. There exists a matrix $A \in \mathbb{R}^{n \times d}$ and a vector $b \in \mathbb{R}^n$ such that, if $m \lesssim d\sqrt{d}$, then the "true" solution $x^* = A^\dagger b$ and the approximation $x' = (DA)^\dagger Db$ have large $\ell_\infty$ distance with constant probability:*

$$\|x' - x^*\|_\infty \gtrsim \frac{1}{\sqrt{d}}\|b\|_2.$$

*Therefore there exists a constant $c$ for which this matrix does not satisfy the generalization guarantee* (6) *with $1 - \frac{c}{d}$ probability.*

*Proof.* We choose the matrix $A$ to be the identity on its top $d$ rows, and $L$ scaled identity matrices $\frac{1}{\sqrt{\alpha d}}I_d$ for the next $dL$ rows, where $L$ satisfies $\frac{1}{d} + \frac{1}{\alpha d}L = 1$ (to normalize each column of $A$), which implies $L = \alpha(d-1)$. Choose some $\beta \in [1, d)$. Set the value of the first $d$ coordinates of vector $b$ to be $\frac{1}{\sqrt{d}}$ and set the value to be $\frac{1}{\sqrt{\beta}}$ for the next $\beta$ coordinates, with the remaining entries all zero. Note that $\|b\|_2 = \sqrt{2}$.

First, we compute $\|Ax - b\|_2^2$. Because $\beta$ is less than $d$, there are two kinds of $x_j$: one involves the following term,

$$\left(\frac{1}{\sqrt{d}}x_j - \frac{1}{\sqrt{d}}\right)^2 + (L-1)\left(\frac{1}{\sqrt{\alpha d}}x_j\right)^2, \tag{15}$$

where the optimal $x_j$ should be set to $1/d$. The other involves the term:

$$\left(\frac{1}{\sqrt{d}}x_j - \frac{1}{\sqrt{d}}\right)^2 + \left(\frac{1}{\sqrt{\alpha d}}x_j - \frac{1}{\sqrt{\beta}}\right)^2 + (L-1)\left(\frac{1}{\sqrt{\alpha d}}x_j\right)^2, \tag{16}$$

where the optimal $x_j$ should be set to $1/d + 1/\sqrt{\alpha\beta d}$. Because we are able to choose $\alpha, \beta$ such that $\alpha\beta \gtrsim d$, then

$$x_j = 1/d + 1/\sqrt{\alpha\beta d} \lesssim 1/d.$$

Second, we compute $\|DAx - Db\|_2^2$. With high probability, there exists a $j$ satisfying Equation (16), but after applying leverage score sampling, the middle term of Equation (16) is removed. Let $p_1 = \frac{1}{d}$ denote the leverage score of each of the top $d$ rows of $A$, and let $p_2 = \frac{1}{\alpha d}$ denote the leverage score of each of the next $Ld$ rows of $A$. We need to discuss the cases $m > d$ and $m \leq d$ separately.

If $m > d$, then the following term involves $x_j$,

$$
\begin{aligned}
&\left(\frac{1}{\sqrt{p_1}}\frac{1}{\sqrt{d}}x_j - \frac{1}{\sqrt{p_1}}\frac{1}{\sqrt{d}}\right)^2 + \frac{m-d}{d} \cdot \left(\frac{1}{\sqrt{p_2}}\frac{1}{\sqrt{\alpha d}}x_j\right)^2 \\
= \ &\frac{1}{p_1}\left(\frac{1}{\sqrt{d}}x_j - \frac{1}{\sqrt{d}}\right)^2 + \frac{m-d}{d} \cdot \frac{1}{p_2}\left(\frac{1}{\sqrt{\alpha d}}x_j\right)^2 \\
= \ &d\left(\left(\frac{1}{\sqrt{d}}x_j - \frac{1}{\sqrt{d}}\right)^2 + \frac{m-d}{d}\alpha\left(\frac{1}{\sqrt{\alpha d}}x_j\right)^2\right).
\end{aligned}
$$

where the optimal $x_j$ should be set to

$$x_j = \frac{1/d}{1/d + (m-d)\alpha/(\alpha d^2)}$$
$$= \frac{1}{1 + (m-d)/d}$$
$$\gtrsim \frac{1}{(m-d)/d}$$
$$\gg \frac{1}{\sqrt{d}}. \hspace{3cm} \text{by } m \ll d\sqrt{d}$$

If $m \leq d$, then the term involving $x_j$ is $(\frac{1}{\sqrt{p_1}} \frac{1}{\sqrt{d}} x_j - \frac{1}{\sqrt{p_1}} \frac{1}{\sqrt{d}})^2$ where the optimal $x_j$ should be set to be $1 \gg 1/\sqrt{d}$.

Third, we need to compute $\|Ax^* - b\|_2^2$ and $\sigma_{\min}(A)$. It is easy to see that $\sigma_{\min}(A)$ because $A$ is an orthonormal matrix. The upper bound for $\|Ax^* - b\|_2^2 = 2$, and the lower bound is also a constant, which can be proved in the following way:

$$\|Ax^* - b\|_2^2 = \sum_{j=1}^{\beta} (15) + \sum_{j=\beta+1}^{d} (16) \geq d(\frac{1}{\sqrt{d}} \frac{1}{d} - \frac{1}{\sqrt{d}})^2 \gtrsim d \cdot \frac{1}{d} = 1.$$

$\blacksquare$

# E  Bounding $\mathbf{E}[\|Z\|_F^2]$

Before getting into the proof details, we define the key property of $S$ being used in the rest of the proofs.

**Definition 24** (All Inner Product Small(AIPS) Property). *For any matrix $S \in \mathbb{R}^{r \times n}$, if for all $i, j \in [n]$ with $i \neq j$ we have*

$$|\langle S_i, S_j \rangle| = O(\sqrt{\log n}/\sqrt{r}),$$

*we say that $S$ satisfies the "AIPS" property.*

**Claim 25.** *If $S \in \mathbb{R}^{r \times n}$ is a subsampled Hadamard transform matrix, then the AIPS property holds with probability at least $1 - 1/\mathrm{poly}(n)$.*

*Proof.* From the structure of $S$, for any $i \neq j$, we have with probability $1 - 1/\mathrm{poly}(n)$ such that $|\langle S_i, S_j \rangle| = O(\sqrt{\log n}/\sqrt{r})$. Applying a union bound over $O(n^2)$ pairs, we obtain that

$$\Pr[\text{ AIPS holds }] \geq 1 - 1/\mathrm{poly}(n).$$

$\blacksquare$

The main idea for bounding $\mathbf{E}[\|Z\|_F^2]$ is to rewrite it as $\mathbf{E}[\|Z\|_F^2] = \mathbf{E}[\|Z\|_F^2 \mid \text{AIPS holds}] + \mathbf{E}[\|Z\|_F^2 \mid \text{AIPS does not hold}]$. Because $\Pr[\text{AIPS does not hold}]$ is at most $1/\mathrm{poly}(n)$, the first term dominates the second term, which means we only need to pay attention to the first term. We repeatedly apply this idea until all the $S$ are removed.

We start by boundinbg $\mathbf{E}[\|Z\|_F^2]$ by squaring $Z_{i_0,j_0}$ and using that $\mathbf{E}[\sigma_i \sigma_j] = 1$ if $i = j$ and $0$ otherwise. Then, we obtain,

$$\mathbf{E}_\sigma [Z_{i_0,j_k}^2] = a_{i_0}^2 b_{j_k}^2 \sum_{i_1, \cdots i_k, j_0, \cdots j_{k-1}} \prod_{c=0}^{k} \langle S_{i_c}, S_{j_c} \rangle^2 \prod_{c=1}^{k} (UU^\top)_{i_{c-1}, j_c}^2. \hspace{1cm} (17)$$

We thus have,

$$\sum_{i_0,j_k,j_k\neq i_k} a_{i_0}^2\langle S_{i_0},S_{j_0}\rangle^2 b_{j_k}^2\langle S_{i_k},S_{j_k}\rangle^2 = a_{j_0}^2\|b\|_2^2 O((\log n)/r) + \|a\|_2^2\|b\|_2^2 O((\log^2 n)/r^2) \overset{\text{def}}{=} C_{j_0},$$

where the first equality is from our conditioning, and the second equality is the definition of $C_{j_0}$. Hence, $\underset{S}{\mathbf{E}}[\|Z\|_F^2]$ is

$$= \sum_{i_1,\cdots i_k,j_0,\cdots j_{k-1}} \prod_{c=1}^{k-1}\langle S_{i_c},S_{j_c}\rangle^2 \prod_{c=1}^{k}(UU^\top)_{j_{c-1},i_c}^2 \cdot \sum_{i_0,j_k,j_k\neq i_k} a_{i_0}^2\langle S_{i_0},S_{j_0}\rangle^2 b_{j_k}^2\langle S_{i_k},S_{j_k}\rangle^2$$

$$= \sum_{i_1,\cdots i_k,j_0,\cdots j_{k-1}} \prod_{c=1}^{k-1}\langle S_{i_c},S_{j_c}\rangle^2 \prod_{c=1}^{k}(UU^\top)_{j_{c-1},i_c}^2 C_{j_0}$$

$$= \sum_{i_1,\cdots i_k,j_0,\cdots j_{k-1}} \langle S_{i_{k-1}},S_{j_{k-1}}\rangle^2 (UU^\top)_{j_{k-1},i_k}^2 \cdot \prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2 \prod_{c=1}^{k-1}(UU^\top)_{j_{c-1},i_c}^2 C_{j_0},$$

where the first equality follows from (17), the second equality by definition of $C_{j_0}$, and the final equality by factoring out $c = k-1$ from one product and $c = k-2$ from the other product.

The way to bound the term $\langle S_{i_{k-1}},S_{j_{k-1}}\rangle$ is by separating the diagonal term where $i_{k-1} = j_{k-1}$ and the non-diagonal term where $i_{k-1} \neq j_{k-1}$. We now use the aforementioned property of $S$, namely, that $\langle S_{i_{k-1}},S_{j_{k-1}}\rangle = 1$, if $i_{k-1} = j_{k-1}$, while for $i_{k-1} \neq j_{k-1}$, we have with probability $1 - 1/\text{poly}(n)$ that $|\langle S_{i_{k-1}},S_{j_{k-1}}\rangle| = O(\sqrt{\log n}/\sqrt{r})$ conditioned on AIPS holding.

Conditioned on AIPS holding, we can recursively reduce the number of terms in the product:

$$\|Z\|_F^2$$

$$= \sum_{i_1,\cdots i_k,j_0,\cdots j_{k-1},i_{k-1}\neq j_{k-1}} O((\log n)/r)\cdot(UU^\top)_{j_{k-1},i_k}^2 \cdot \prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2 \prod_{c=1}^{k-1}(UU^\top)_{j_{c-1},i_c}^2 C_{j_0}$$

$$+ \sum_{i_1,\cdots i_k,j_0,\cdots j_{k-1},i_{k-1}=j_{k-1}} 1\cdot(UU^\top)_{j_{k-1},i_k}^2 \cdot \prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2 \prod_{c=1}^{k-1}(UU^\top)_{j_{c-1},i_c}^2 C_{j_0}$$

$$\leq \sum_{i_1,\cdots i_k,j_0,\cdots j_{k-1}} O((\log n)/r)\cdot(UU^\top)_{j_{k-1},i_k}^2 \cdot \prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2 \prod_{c=1}^{k-1}(UU^\top)_{j_{c-1},i_c}^2 C_{j_0}$$

$$+ \sum_{i_1,\cdots i_k,j_0,\cdots j_{k-1},i_{k-1}=j_{k-1}} 1\cdot(UU^\top)_{j_{k-1},i_k}^2 \cdot \prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2 \prod_{c=1}^{k-1}(UU^\top)_{j_{c-1},i_c}^2 C_{j_0},$$

where the first equality follows from the property just mentioned, and the inequality follows by including back the tuples of indices for which $i_{k-1} = j_{k-1}$, using that each summand is non-negative.

Our next step will be to bound the term $(UU^\top)^2_{j_{k-1},i_k}$. We have, $\|Z\|_F^2$ is

$$\leq \sum_{i_k,j_{k-1}} (UU^\top)^2_{i_k,j_{k-1}} \sum_{\substack{i_1,\cdots,i_{k-1}\\ j_0,\cdots,j_{k-2}}} O((\log n)/r)$$

$$\cdot \prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2 \prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c}C_{j_0}$$

$$+ \sum_{\substack{i_1,\cdots,i_k\\ j_0,\cdots,j_{k-1}\\ i_{k-1}=j_{k-1}}} 1\cdot(UU^\top)^2_{j_{k-1},i_k}\prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2\prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c}C_{j_0}$$

$$= O(d(\log n)/r)\underbrace{\sum_{\substack{i_1,\cdots,i_{k-1}\\ j_0,\cdots,j_{k-2}}}\prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2\prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c}C_{j_0}}_{A}$$

$$+ \underbrace{\sum_{\substack{i_1,\cdots,i_k\\ j_0,\cdots,j_{k-1}\\ i_{k-1}=j_{k-1}}} 1\cdot(UU^\top)^2_{j_{k-1},i_k}\prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2\prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c}C_{j_0}}_{B},$$

where the equality uses that $\sum_{i_k,j_{k-1}}(UU^\top)^2_{i_k,j_{k-1}} = \|UU^\top\|_F^2 = d$. We first upper bound term $B$:

$$= \sum_{\substack{i_1,\cdots,i_k\\ j_0,\cdots,j_{k-1}\\ i_{k-1}=j_{k-1}}} 1\cdot(UU^\top)^2_{j_{k-1},i_k}\prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2\prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c}C_{j_0}$$

$$= \sum_{\substack{i_1,\cdots,i_{k-1}\\ j_0,j_1,\cdots,j_{k-1}\\ i_{k-1}=j_{k-1}}} C_{j_0}\prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2\prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c}\sum_{i_k}(UU^\top)^2_{j_{k-1},i_k}$$

$$= \sum_{\substack{i_1,\cdots,i_{k-1}\\ j_0,j_1,\cdots,j_{k-1}\\ i_{k-1}=j_{k-1}}} C_{j_0}\prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2\prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c}|e_{j_{k-1}}UU^\top|^2$$

$$\leq \sum_{\substack{i_1,\cdots,i_{k-1}\\ j_0,j_1,\cdots,j_{k-1}\\ i_{k-1}=j_{k-1}}} C_{j_0}\prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2\prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c}1$$

$$= \sum_{\substack{i_1,\cdots,i_{k-1}\\ j_0,j_1,\cdots,j_{k-2}}} C_{j_0}\prod_{c=1}^{k-2}\langle S_{i_c},S_{j_c}\rangle^2\prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c},$$

where the first equality is the definition of $B$, the second equality follows by separating out the

index $i_k$, the third equality uses that $\sum_{i_k}(UU^\top)^2_{j_{k-1},i_k} = \|e_{j_{k-1}}UU^\top\|^2_2$, that is, the squared norm of the $j_{k-1}$-th row of $UU^\top$, the inequality follows since all rows of a projection matrix $UU^\top$ have norm at most 1, and the final equality uses that $j_{k-1}$ no longer appears in the expression.

We now merge our bounds for the terms $A$ and $B$ in the following way:

$$
\begin{aligned}
&\|Z\|^2_F \\
\leq\ & A + B \\
\leq\ & O(d(\log n)/r) \sum_{\substack{i_1,\cdots,i_{k-1} \\ j_0,\cdots,j_{k-2}}} \prod_{c=1}^{k-2}\langle S_{i_c}, S_{j_c}\rangle^2 \prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c} C_{j_0} \\
&+ \sum_{\substack{i_1,\cdots,i_{k-1} \\ j_0,j_1,\cdots,j_{k-2}}} C_{j_0} \prod_{c=1}^{k-2}\langle S_{i_c}, S_{j_c}\rangle^2 \prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c} \\
=\ & (O(d(\log n)/r)+1) \sum_{\substack{i_1,\cdots,i_{k-1} \\ j_0,\cdots,j_{k-2}}} \prod_{c=1}^{k-2}\langle S_{i_c}, S_{j_c}\rangle^2 \prod_{c=1}^{k-1}(UU^\top)^2_{j_{c-1},i_c} C_{j_0} \\
\leq\ & \cdots \\
\leq\ & (O(d(\log n)/r)+1)^2 \sum_{\substack{i_1,\cdots,i_{k-2} \\ j_0,\cdots,j_{k-3}}} \prod_{c=1}^{k-3}\langle S_{i_c}, S_{j_c}\rangle^2 \prod_{c=1}^{k-2}(UU^\top)^2_{j_{c-1},i_c} C_{j_0} \\
\leq\ & \cdots \\
\leq\ & (O(d(\log n)/r)+1)^{k-1} \sum_{i_1,j_0} \prod_{c=1}^{1}(UU^\top)^2_{j_{c-1},i_c} C_{j_0} \\
\leq\ & (O(d(\log n)/r)+1)^{k-1}\left(d\|b\|^2_2(\log^2 n)/r^2 + \|b\|^2_2(\log n)/r\right),
\end{aligned}
$$

where the first two inequalities and first equality are by definition of $A$ and $B$ above. The first inequality follows by induction, since at this point we have replaced $k$ with $k-1$, and can repeat the argument, incurring another multiplicative factor of $O(d(\log n)/r)+1$. Repeating the induction in this way we arrive at the last inequality. Finally, the last inequality follows by plugging in the definition of $C_{j_0}$, using that $\sum_{i_1,j_0}(UU^\top)^2_{j_0,i_1} = d$, and

$$
\sum_{j_0,i_1}(UU^\top)^2_{j_0,i_1}a^2_{j_0} = \sum_{j_0}a^2_{j_0}\sum_{i_1}(UU^\top)^2_{j_0,i_1} = \sum_{j_0}a^2_{j_0}\|e_{j_0}UU^\top\|^2_2 \leq 1,
$$

where the inequality follows since each row of $UU^\top$ has norm at most 1, and $a$ is a unit vector. The final result is that

$$
\|Z\|^2_F \leq (O(d(\log n)/r)+1)^{k-1}\left(d\|b\|^2_2(\log^2 n)/r^2 + \|b\|^2_2(\log n)/r\right).
$$