

Information-Theoretic Methods in Data Science

Information-theoretic bounds on sketching

Mert Pilanci

Department of Electrical Engineering
Stanford University

Contents

1	Information-theoretic bounds on sketching	<i>page</i> 3
1.1	Introduction	3
1.2	Types of randomized sketches	4
1.2.1	Gaussian sketches	5
1.2.2	Sub-Gaussian sketches	5
1.2.3	Randomized orthonormal systems (ROS)	5
1.2.4	Sketches based on random row sampling:	6
1.2.5	Graph sparsification via sub-sampling	7
1.2.6	Sparse sketches based on hashing	7
1.3	Background on convex analysis and optimization	8
1.4	Sketching upper-bounds for regression problems	9
1.4.1	Overdetermined case ($n > d$)	9
1.4.2	Gaussian complexity	10
1.4.3	Underdetermined case ($n \leq d$)	12
1.5	Information-theoretic lower-bounds	13
1.5.1	Statistical upper and lower bounds	13
1.5.2	Fano's inequality	15
1.5.3	Metric entropy	15
1.5.4	Minimax risk	16
1.5.5	Reduction to hypothesis testing	16
1.5.6	Implications of the information-theoretic lower bound	17
1.5.7	Iterative sketching	20
1.6	Nonparametric problems	21
1.6.1	Nonparametric regression	21
1.6.2	Sketching kernels	23
1.7	Extensions: Privacy and communication complexity	24
1.7.1	Privacy and information-theoretic bounds	24
1.7.2	Mutual information privacy	24
1.7.3	Optimization based privacy attacks	25
1.7.4	Communication complexity space lower bounds	25
1.8	Numerical experiments	26
1.9	Summary	27
1.10	Proof of Theorem 1.6	28

1.11 Proof of Lemma 1.5

29

References

31

1 Information-theoretic bounds on sketching

1.1 Introduction

In the recent years we have witnessed an unprecedented increase in the amount of available data in a wide variety of fields. Approximate computation methods with provable performance guarantees are becoming important and relevant tools in practice to attack larger scale problems. The term *sketching* is used for randomized algorithms designed to reduce data dimensionality in computationally intensive tasks. Sketching can often provide better space, time and communication complexity trade offs by sacrificing minimal accuracy. In large scale data science problems, sketching allows us to leverage limited computational resources such as memory, time and bandwidth, and also explore favorable trade-offs between accuracy and computational complexity.

Random projections are widely used instances of sketching, and gathered substantial attention in the literature, especially very recently in machine learning, signal processing and theoretical computer science communities [1, 2, 3, 4, 5, 6]. Other popular sketching techniques include leverage score sampling, graph sparsification, core-sets and randomized matrix factorizations. In this chapter we overview sketching methods, develop lower bounds using information-theoretic techniques and present upper-bounds on their performance. In the next section we begin by introducing commonly used sketching methods.

This chapter will be focused on the role of information theory in sketching methods for solving large scale statistical estimation and optimization problems, and investigate fundamental lower-bounds on their performance. By exploring these lower-bounds, we obtain interesting trade-offs in computation and accuracy. Moreover, we may hope to obtain improved sketching constructions by understanding their information-theoretic properties. The lower-bounding techniques employed here parallel the information-theoretic techniques used in statistical minimax theory [7, 8]. We apply Fano's inequality and packing constructions to understand fundamental lower-bounds on the accuracy of sketching.

Randomness and sketching also have applications in privacy preserving queries [9, 10]. Privacy has become an important concern in the age of information where breach of sensitive data is frequent. We will illustrate that randomized sketching offers a computationally simple and effective mechanism to preserve privacy in optimization and machine learning.

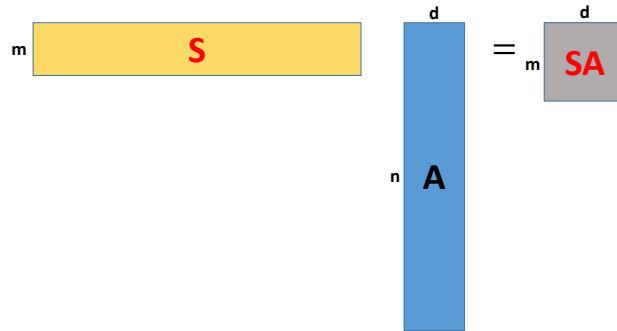


Figure 1.1 Sketching a tall matrix A . The smaller matrix $SA \in \mathbb{R}^{m \times d}$ is a compressed version of the original data $A \in \mathbb{R}^{n \times d}$

We start with an overview of different constructions of sketching matrices in Section 1.2. Then in Section 1.3 we present upper-bounds on the performance of sketching from an optimization viewpoint. To be able to analyze upper-bounds, we introduce the notion of *localized Gaussian complexity*, which also plays an important role in the characterization of minimax statistical bounds. In section 1.4 we will discuss information-theoretic lower-bounds on sketching from a statistical perspective. In Section 1.5 we turn to nonparametric problems and information-theoretic lower-bounds. Finally, in Section 1.6 we present privacy preserving properties of sketching using a mutual information characterization.

1.2 Types of randomized sketches

In this section we describe popular constructions of sketching matrices. Given a sketching matrix S , we use $\{s_i\}_{i=1}^m$ to denote the collection of its n -dimensional rows. Here we consider sketches which are zero-mean, and are normalized, i.e., they satisfy the following two conditions

$$(a) \quad \mathbb{E} S^T S = I_{d \times d} \quad (1.1)$$

$$(b) \quad \mathbb{E} S = 0_{n \times d} \quad (1.2)$$

These reasoning behind the above conditions will be more clear when applied to sketching optimization problems involving data matrices.

A very typical use of sketching is to obtain compressed versions of a large data matrix A . We obtain the matrix $SA \in \mathbb{R}^{m \times d}$ using simple matrix multiplication. See Figure 1.1 for an illustration. As we will see in a variety of examples, random matrices preserve most of the information in the matrix A .

1.2.1 Gaussian sketches

The most classical sketch is based on a random matrix $S \in \mathbb{R}^{m \times n}$ with i.i.d. standard Gaussian entries. Suppose that we generate a random matrix $S \in \mathbb{R}^{m \times n}$ with entries drawn from i.i.d. zero-mean Gaussian random variables with variance $\frac{1}{m}$. Note that we have $\mathbb{E}S = 0_{m \times d}$ and also $\mathbb{E}S^T S = \sum_{i=1}^m \mathbb{E}s_i s_i^T = \sum_{i=1}^m I_d \frac{1}{m} = I_d$. Analyzing the Gaussian sketches is considerably easier than others, due to special properties of the Gaussian distribution such as rotation invariance. However, Gaussian sketches may not be the most computationally efficient choice for many data matrices as we will see in the sequel.

1.2.2 Sub-Gaussian sketches

A generalization of the previous construction is a random sketch with row drawn from i.i.d. sub-Gaussian random variables. In particular, a zero-mean random vector $s \in \mathbb{R}^n$ is 1-sub-Gaussian if for any $u \in \mathbb{R}^n$, we have

$$\mathbb{P}[\langle s, u \rangle \geq \varepsilon \|u\|_2] \leq e^{-\varepsilon^2/2} \quad \text{for all } \varepsilon \geq 0. \quad (1.3)$$

For instance, a vector with i.i.d. $N(0, 1)$ entries is 1-sub-Gaussian, as is a vector with i.i.d. Rademacher entries (uniformly distributed over $\{-1, +1\}$). In many models of computation, multiplying numbers with random signs is simpler than multiplying with Gaussian variables, and only costs an addition operation. Note that multiplying with -1 only amounts to flipping the sign bit in the signed number representation of the number in the binary system. In modern computers, the difference between addition and multiplication is often inappreciable. However, the real disadvantage of sub-Gaussian and Gaussian sketches is that, they require matrix-vector multiplications with unstructured and dense random matrices. In particular, given a data matrix $A \in \mathbb{R}^{n \times d}$, computing its sketched version SA requires $\mathcal{O}(mnd)$ basic operations using classical matrix multiplication algorithms, in general.

1.2.3 Randomized orthonormal systems (ROS)

The second type of randomized sketch we consider is *randomized orthonormal system* (ROS), for which matrix multiplication can be performed much more efficiently.

In order to define a ROS sketch, we first let $H \in \mathbb{R}^{n \times n}$ be an orthonormal matrix with entries $H_{ij} \in [-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}]$. Standard classes of such matrices are the Hadamard or Fourier bases, for which matrix-vector multiplication can be performed in $\mathcal{O}(n \log n)$ time via the fast Hadamard or Fourier transforms, respectively. For example, an $n \times n$ Hadamard matrix $H = H_n$ can be recursively

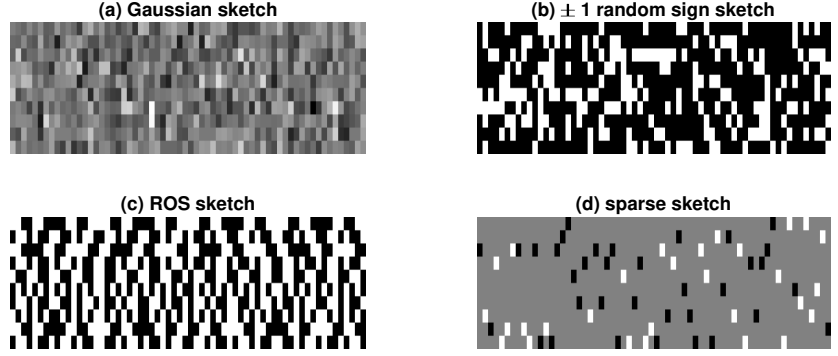


Figure 1.2 Different types of sketching matrices: (a) Gaussian sketch, (b) ± 1 random sign sketch, (c) Randomized Orthogonal System sketch, (d) sparse sketch

constructed as follows:

$$H_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad H_4 = \frac{1}{\sqrt{2}} \begin{bmatrix} H_2 & H_2 \\ H_2 & -H_2 \end{bmatrix} \quad H_{2^t} = \underbrace{H_2 \otimes H_2 \otimes \cdots \otimes H_2}_{\text{Kronecker product } t \text{ times}}$$

Based on any such matrix, a sketching matrix $S \in \mathbb{R}^{m \times n}$ from a ROS ensemble is obtained by sampling i.i.d. rows of the form

$$s^T = \sqrt{n} e_j^T H D \quad \text{with probability } 1/n \text{ for } j = 1, \dots, n,$$

where the random vector $e_j \in \mathbb{R}^n$ is chosen uniformly at random from the set of all n canonical basis vectors, and $D = \text{diag}(\nu)$ is a diagonal matrix of i.i.d. Rademacher variables $\nu \in \{-1, +1\}^n$. Alternatively, the rows of the ROS sketch can be sampled without replacement and one can obtain similar guarantees to sampling with replacement. Given a fast routine for matrix-vector multiplication, ROS sketch SA of the data $A \in \mathbb{R}^{n \times d}$ can be formed in $\mathcal{O}(nd \log m)$ time (for instance, see [11]).

1.2.4 Sketches based on random row sampling:

Given a probability distribution $\{p_j\}_{j=1}^n$ over $[n] = \{1, \dots, n\}$, another choice of sketch is to randomly sample the rows of the extended data matrix A a total of m times with replacement from the given probability distribution. Thus, the rows of S are independent and take on the values

$$s^T = \frac{e_j^T}{\sqrt{p_j}} \quad \text{with probability } p_j \text{ for } j = 1, \dots, n$$

where $e_j \in \mathbb{R}^n$ is the j^{th} canonical basis vector. Different choices of the weights $\{p_j\}_{j=1}^n$ are possible, including those based on the *leverage values* of A . Leverage

values are defined as

$$p_j := \frac{\|u_j\|_2^2}{\sum_{i=1}^n \|u_i\|_2^2}$$

where u_1, u_2, \dots, u_n are the rows of $U \in \mathbb{R}^{n \times d}$ which is the matrix of left singular vectors of A . Leverage values can be obtained using a Singular Value Decomposition $A = U\Sigma V^T$. Moreover, there also exists faster randomized algorithms to approximate the leverage scores (e.g., see [12]). In our analysis of lower bounds to follow, we assume that the weights are α -balanced, meaning that

$$\max_{j=1, \dots, n} p_j \leq \frac{\alpha}{n} \quad (1.4)$$

for some constant α independent of n .

1.2.5 Graph sparsification via sub-sampling

Let $G = (V, E)$ be a weighed, undirected graph with d nodes and n edges, where V and E are the set of nodes and edges respectively. Let $A \in \mathbb{R}^{n \times d}$ be the node-edge incidence matrix of the graph A . Suppose we randomly sample the edges in the graph a total of m times with replacement from a given probability distribution over the edges. The obtained graph is a weighted sub-graph of the original, whose incidence matrix is SA . Similar to the row sampling sketch, the sketch can be written as

$$s^T = \frac{e_j^T}{\sqrt{p_j}} \quad \text{with probability } p_j \text{ for } j = 1, \dots, n.$$

We note that row and graph sub sampling sketches satisfy the condition (1.1). However, they do not satisfy the condition (1.2). In many computational problems on graphs, sparsifying the graph has computational advantages. Notable examples of such problems are solving Laplacian linear systems and graph partitioning where sparsification can be used. We refer the reader to Spielman and Srivastava [13] for details.

1.2.6 Sparse sketches based on hashing

In many applications, the data matrices contain very few non-zero entries. For sparse data matrices, special constructions of the sketching matrices yield greatly improved performance. Here we describe the count-sketch construction from [14, 15]. Let $h : [n] \rightarrow [m]$ be a hash functions from a pairwise independent family¹. The entry S_{ij} of the sketch matrix is given by σ_j if $i = h(j)$, and otherwise it is zero, where $\sigma \in \{-1, +1\}^m$ is a random vector containing 4-wise independent variables. Therefore, the j^{th} column of S is nonzero only in row indexed by $h(j)$. We refer the reader to [14, 15] for the details. An example realization of the

¹ A hash function is from a pairwise independent family if $\mathbb{P}[h(j) = i, h(k) = l] = \frac{1}{m^2}$ and $\mathbb{P}[h(j) = i] = \frac{1}{m}$ for all i, j, k, l .

sparse sketch is given below, where each column contains a single non-zero entry which is uniformly random sign ± 1 at a uniformly random index.

$$\begin{bmatrix} 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & -1 \end{bmatrix}.$$

Figure 1.2 shows examples of different randomized sketching matrices $S \in \mathbb{R}^{m \times n}$ where $m = 64, n = 1024$ which are drawn randomly. We refer readers to Nelson and Nguyen [16] for details on sparse sketches.

1.3 Background on convex analysis and optimization

In this section, we first briefly review relevant concepts from convex analysis and optimization. A set $\mathcal{C} \subseteq \mathbb{R}^d$ is convex if for any $x, y \in \mathcal{C}$

$$tx + (1 - t)y \in \mathcal{C} \text{ for all } t \in [0, 1].$$

Let X be a convex set. A function $f : X \rightarrow \mathbb{R}$ is convex if for any $x, y \in X$

$$f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y) \text{ for all } t \in [0, 1].$$

Given a matrix $A \in \mathbb{R}^{n \times d}$, we define the linear transform of the convex set \mathcal{C} as $A\mathcal{C} = \{Ax \mid x \in \mathcal{C}\}$. It can be shown that $A\mathcal{C}$ is convex if \mathcal{C} is convex.

A convex optimization problem is a minimization problem of the form

$$\min_{x \in \mathcal{C}} f(x) \tag{1.5}$$

where $f(x)$ is a convex function and \mathcal{C} is a convex set. In order to characterize optimality of solutions, we will define the tangent cone of \mathcal{C} at a fixed vector x^* as follows

$$\mathcal{T}_{\mathcal{C}}(x^*) = \{t(x - x^*) \mid t \geq 0 \text{ and } x \in \mathcal{C}\}. \tag{1.6}$$

Figures 1.3 and 1.4 illustrate two examples of tangent cones of a polyhedral convex set in \mathbb{R}^2 .² A first order characterization of optimality in the convex optimization problem (1.5) is given by the tangent cone. If a vector x^* is optimal in (1.5), it holds that

$$z^T \nabla f(x^*) \geq 0, \forall z \in \mathcal{T}_{\mathcal{C}}(x^*). \tag{1.7}$$

We refer the reader to Hiriart-Urruty and Lemaréchal [17] for details on convex analysis, and Boyd and Vandenberghe [18] for an in-depth discussion of convex optimization problems and applications.

² Note that the tangent cones extend towards infinity in certain directions, whereas the shaded regions in Figures 1.3 and 1.4 are compact for illustration.

1.4 Sketching upper-bounds for regression problems

Now we consider an instance of a convex optimization problem. Consider the Least Squares optimization

$$x^* = \arg \min_{x \in \mathcal{C}} \underbrace{\|Ax - b\|_2}_{f(x)}, \quad (1.8)$$

where $A \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$ are the input data and $\mathcal{C} \subseteq \mathbb{R}^d$ is a closed and convex constraint set. In statistical and signal processing applications, it is typical to use the constraint set to impose structure on the obtained solution x . Important examples of the convex constraint \mathcal{C} include the nonnegative orthant, ℓ_1 -ball for promoting sparsity, and ℓ_∞ -ball as a relaxation to the combinatorial set $\{0, 1\}^d$.

In the unconstrained case when $\mathcal{C} = \mathbb{R}^d$, a closed form solution exists for the solution, which is given by $x^* = (A^T A)^{-1} A^T b$. However, forming the Gram matrix $A^T A$ and inverting using direct methods such as QR decomposition, or the Singular Value Decomposition typically requires $O(nd^2) + O(nd \min(n, d))$ operations respectively. Faster iterative algorithms such as the conjugate gradient (CG) method can be used to obtain an approximate solution in $O(nd\kappa(A))$ time, where $\kappa(A)$ is the condition number of the data matrix A . Using sketching methods, it is possible to obtain even faster approximate solutions as we will discuss in the sequel.

In the constrained case, a variety of efficient iterative algorithms have been developed in the last couple decades to obtain the solution, such as proximal and projected gradient methods, their accelerated variants, and barrier based second-order methods. Sketching can also be used to improve the run-time of these methods.

1.4.1 Overdetermined case ($n > d$)

In many of the applications, the number of observations n , exceed the number of unknowns d , which gives rise to the tall $n \times d$ matrix A . In machine learning, it is very common to encounter datasets where n is very large, and d is of moderate size. Suppose that we first compute the sketched data matrices SA and Sb from the original data A and b , then consider the following approximation to the above optimization problem

$$\hat{x} = \arg \min_{x \in \mathcal{C}} \|SAx - Sb\|_2^2. \quad (1.9)$$

After applying the sketch to the data matrices, the sketched problem has dimensions $m \times d$, which is lower than the original dimensions when $m < n$. Note that the objective in the above problem (1.9) can be seen as an unbiased approximation of the original objective function (1.8), since it holds that

$$\mathbb{E} \|SAx - Sb\|_2^2 = \|Ax - b\|_2^2.$$

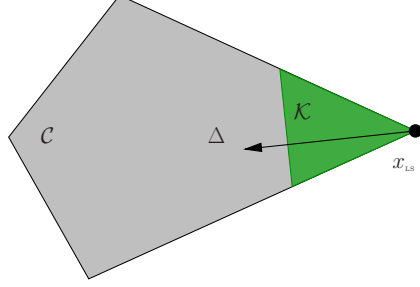


Figure 1.3 (a) A narrow tangent cone where the Gaussian complexity is small

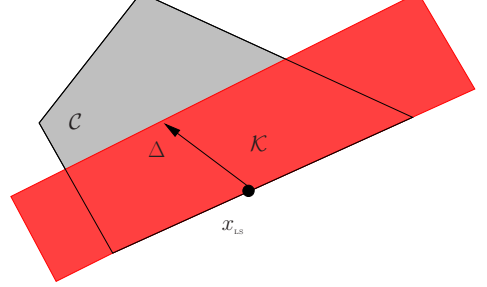


Figure 1.4 (b) A wide tangent cone where the Gaussian complexity is large

for any fixed choice of A, x and b . This is a consequence of the condition (1.1), which is satisfied by all of the sketching matrices considered in Section 1.2.

1.4.2 Gaussian complexity

Gaussian complexity plays an important role in statistics, empirical process theory, compressed sensing and the theory of Banach spaces [19, 20, 21]. Here we consider a localized version of the Gaussian complexity, which is defined as follows

$$\mathcal{W}_t(\mathcal{C}) := \mathbb{E}_g \left[\sup_{\substack{\theta \in \mathcal{C} \\ \|\theta\|_2 \leq t}} |\langle g, \theta \rangle| \right], \quad (1.10)$$

where g is a random vector with i.i.d. standard Gaussian entries, i.e., $g \sim \mathcal{N}(0, I_n)$. The parameter $t > 0$ controls the radius at which the random deviations are localized. For a finite value of t , the supremum in (1.10) is always achieved since the constraint set is compact.

Analyzing the sketched optimization problem requires us to control the random deviations constrained to the set of possible descent directions $\{x - x^* \mid x \in \mathcal{C}\}$. We now define a transformed tangent cone at x^* as follows

$$\mathcal{K} = \{tA(x - x^*) \mid t \geq 0 \text{ and } x \in \mathcal{C}\},$$

which can be alternatively defined as $A\mathcal{T}_{\mathcal{C}}(x^*)$ using the definition given in (1.6). The next theorem provides an upper-bound on the performance of the sketching method for constrained optimization based on localized Gaussian complexity.

THEOREM 1.1 *Let S be a Gaussian sketch, and let \hat{x} be the solution of (1.9). Suppose that $m \geq \frac{c_0 \mathcal{W}_1(\mathcal{K})^2}{\epsilon^2}$, where c_0 is a universal constant, then it holds that*

$$\frac{\|A(\hat{x} - x^*)\|_2}{f(x^*)} \leq \epsilon,$$

and consequently we have

$$f(x^*) \leq f(\hat{x}) \leq f(x^*)(1 + \epsilon). \quad (1.11)$$

As predicted by the theorem, the approximation ratio improves as the sketch dimension m increases, and converges to one as $m \rightarrow \infty$. However, we are often interested in the rate of convergence of the approximation ratio. Theorem 1.1 characterizes this rate by relating the geometry of the constraint set to the accuracy of the sketching method (1.9). As an illustration, Figures 1.2 and 1.3 depict a narrow and wider tangent cone. The proof of Theorem 1.1 combines the convex optimality condition involving the tangent cone in (1.7) with results on empirical processes, and can be found in Pilanci and Wainwright [22]. An important feature of Theorem 1.1 is that the approximation quality is relative to the optimal value $f(x^*)$. This is advantageous when $f(x^*)$ is small, e.g., the optimal value can be zero in noise-less signal recovery problems. However, in problems where the signal-to-noise ratio is low, $f(x^*)$ can be large, and hence negatively effect the approximation quality. We illustrate the implications of Theorem 1.1 on some concrete examples in the sequel.

Example 1: Unconstrained least squares

For unconstrained problems, we have $\mathcal{K} = \text{range}(A)$, i.e., the tangent cone is equal to the range space of the data matrix A . In order to apply Theorem 1.1, we need the following lemma about the Gaussian complexity of a subspace.

LEMMA 1.2 *Let Q be a subspace of dimension q . The Gaussian complexity of Q satisfies*

$$\mathcal{W}_t(Q) \leq t\sqrt{q}.$$

Proof Let U be an orthonormal basis for the subspace Q . We have the following representation $L = \{Ux \mid x \in \mathbb{R}^q\}$. Consequently the Gaussian complexity $\mathcal{W}_1(Q)$ can be written as

$$\begin{aligned} \mathbb{E}_g \left[\sup_{\substack{x \\ \|Ux\|_2 \leq t}} \langle g, Ux \rangle \right] &= \mathbb{E}_g \left[\sup_{\substack{x \\ \|x\|_2 \leq t}} \langle U^T g, x \rangle \right] = t \mathbb{E}_g \|U^T g\|_2 \leq t \sqrt{\mathbb{E} \text{tr} U U^T g g^T} \\ &= t \sqrt{\text{tr} U^T U} \\ &= t\sqrt{q}. \end{aligned}$$

Where the inequality follows from Jensen's inequality and concavity of the square root, and first and fifth equality follows since $U^T U = I_q$. Therefore, the Gaussian complexity of the range of A for $t = 1$ satisfies

$$\mathcal{W}_1(\text{range}(A)) \leq \sqrt{\text{rank}(A)}.$$

□

Setting the dimension of the sketch $m \geq c_0 \text{rank}(A)/\epsilon^2$ suffices to obtain an ϵ approximate solution in the sense of (1.11). We note that $\text{rank}(A)$ might not be available a priori, but the upper bound $\text{rank}(A) \leq d$ may be useful when $n \gg d$.

Example 2: ℓ_1 constrained least squares

For ℓ_1 -norm constrained problems we have $\mathcal{C} = \{x \mid \|x\|_1 \leq r\}$ for some radius parameter r . The tangent cone \mathcal{K} at the optimal point x^* depends on the support of x^{*3} , and hence its cardinality $\|x^*\|_0$. In [23] it is shown that the localized Gaussian complexity satisfies

$$\mathcal{W}_1(\mathcal{K}) \leq c_1 \frac{\gamma_+(A)}{\gamma_-(A)} \sqrt{\|x^*\|_0 \log d},$$

where c_1 is a universal constant and γ_k are β_k are ℓ_1 -restricted maximum and minimum eigenvalues defined as follows

$$\gamma_k := \max_{\substack{\|z\|_2=1 \\ \|z\|_1 \leq \sqrt{k}}} \|Az\|_2^2 \quad \text{and} \quad \beta_k := \min_{\substack{\|z\|_2=1 \\ \|z\|_1 \leq \sqrt{k}}} \|Az\|_2^2.$$

As a result, we conclude that for ℓ_1 constrained problems, the sketch dimension can be substantially smaller when ℓ_1 constrained eigenvalues are well behaved.

1.4.3 Underdetermined case ($n \leq d$)

In many applications the dimension of the data vectors may be larger than the sample size. In these situations, it makes sense to reduce the dimensionality by applying the sketch on the right, i.e., AS^T and solve

$$\arg \min_{\substack{z \in \mathbb{R}^m \\ S^T z \in \mathcal{C}}} \|(AS^T z - b)\|_2. \quad (1.12)$$

Note that the vector $z \in \mathbb{R}^m$ is of smaller dimension than the original variable $x \in \mathbb{R}^d$. After solving the reduced dimensional problem and obtaining its optimal solution z^* , the final estimate for the original variable x can be taken as $\hat{x} = S^T z^*$. We will investigate this approach in Section 1.5 in non-parametric statistical estimation problems and present concrete theoretical guarantees.

It is instructive to note that, in the special case where we have ℓ_2 regularization and $\mathcal{C} = \mathbb{R}^d$, we can easily transform the underdetermined least squares problem into an overdetermined one using convex duality, or matrix inversion lemma. We first write the sketched problem (1.12) as the constrained convex program

$$\min_{\substack{z \in \mathbb{R}^m, y \in \mathbb{R}^n \\ y = AS^T z}} \frac{1}{2} \|y - b\|_2^2 + \rho \|z\|_2^2,$$

and form the convex dual. It can be shown that strong duality holds, and consequently primal and dual programs can be stated as follows

$$\min_{z \in \mathbb{R}^m} \frac{1}{2} \|AS^T z - b\|_2^2 + \rho \|z\|_2^2 = \max_{x \in \mathbb{R}^d} -\frac{1}{4\rho} \|SA^T x\|_2^2 - \frac{1}{2} \|x\|_2^2 + x^T b,$$

where the primal and dual solutions satisfy $z^* = \frac{1}{2\rho} SA^T x^*$ at the optimum [18].

³ The term support refers to the set of indices where the solution has a non-zero value

Therefore the sketching matrix applied from the right, AS^T , corresponds to a sketch applied on the left, SA^T , in the dual problem which parallels (1.9). This observation can be used to derive approximation results on the dual program. We refer the reader to [22] for an application in support vector machine classification where $b = 0_n$.

1.5 Information-theoretic lower-bounds

1.5.1 Statistical upper and lower bounds

In order to develop information-theoretic lower-bounds, we consider a statistical observation model for the constrained regression problem. Consider the following model

$$b = Ax^\dagger + w, \quad \text{where } w \sim \mathcal{N}(0, \sigma^2 I_n), \text{ and } x^\dagger \in \mathcal{C}_0, \quad (1.13)$$

where x^\dagger is the unknown vector to be estimated and w is an i.i.d. noise vector whose entries are distributed as $\mathcal{N}(0, \sigma^2)$. In this section we will focus on the observation model (1.13) and present a lower bound on all estimators which use the sketched data (SA, Sb) to form an estimate \hat{x} .

We assume that the unknown vector x^\dagger belongs to some set $\mathcal{C}_0 \subseteq \mathcal{C}$ that is star-shaped around zero⁴. In many cases of interest we have $\mathcal{C} = \mathcal{C}_0$, i.e., when the set \mathcal{C} is convex and simple to describe. In this case, the constrained least-squares estimate x^* from equation (1.8) corresponds to the constrained maximum-likelihood estimator for estimating the unknown regression vector x^\dagger under the Gaussian observation model (1.13). However \mathcal{C}_0 may not be computationally tractable as an optimization constraint set, such as a non-convex set and we can consider a set \mathcal{C} which is a convex relaxation⁵ of this set, such that $\mathcal{C} \subset \mathcal{C}_0$. An important example is the set of s sparse and bounded vectors given by $\mathcal{C}_0 = \{x : \|x\|_0 \leq s, \|x\|_\infty \leq 1\}$, which has combinatorially many elements. The well-known ℓ_1 relaxation given by $\mathcal{C} = \{x : \|x\|_1 \leq \sqrt{s}, \|x\|_\infty \leq 1\}$ satisfies $\mathcal{C} \subset \mathcal{C}_0$, which follows from Cauchy-Schwartz inequality, and is widely used [24, 25] to find sparse solutions.

We now present a theoretical result on the statistical performance of the original constrained least squares estimator in (1.8)

THEOREM 1.3 *Let \mathcal{C} be any set that contains the true parameter x^\dagger . Then the constrained estimator x^* in (1.8) under the observation model (1.13) has mean-squared error upper bound as*

$$\mathbb{E}_w \left[\frac{1}{n} \|A(x^* - x^\dagger)\|_2^2 \right] \leq c_2 \left(\delta^*(n)^2 + \frac{\sigma^2}{n} \right),$$

⁴ Star-shaped assumption means that for any $x \in \mathcal{C}_0$ and scalar $t \in [0, 1]$, the point tx also belongs to \mathcal{C}_0 .

⁵ We may also consider an approximation of \mathcal{C}_0 which doesn't necessarily satisfy $\mathcal{C} \subset \mathcal{C}_0$. An example is the ℓ_1 and ℓ_0 unit balls.

where $\delta^*(n)$ is the critical radius, equal to the smallest positive solution $\delta > 0$ to the inequality

$$\frac{\mathcal{W}_\delta(\mathcal{C})}{\delta\sqrt{n}} \leq \frac{\delta}{\sigma}. \quad (1.14)$$

Please see [23, 20] for a proof of this theorem. This result provides a baseline against which to compare the statistical recovery performance of the randomized sketching method. In particular, an important goal is characterizing the minimal projection dimension m that will enable us to find an estimate \hat{x} with the error guarantee

$$\frac{1}{n} \|A(\hat{x} - x^\dagger)\|_2^2 \approx \frac{1}{n} \|A(x^* - x^\dagger)\|_2^2,$$

in a computationally simpler manner using the compressed data SA, Sb .

An application of Theorem 1.1 will yield that the sketched solution \hat{x} in (1.9) using the choice of sketch dimension $m = \frac{c_0 \mathcal{W}_1(\mathcal{K})^2}{\epsilon^2}$, satisfies the bound

$$\|A(\hat{x} - x^*)\| \leq \epsilon \|Ax^* - b\|_2,$$

where $\|Ax^* - b\|_2 = f(x^*)$ is the optimal value of the optimization problem (1.8). However, under the model (1.13) we have

$$\|Ax^* - b\|_2 = \|A(x^* - x^\dagger) - w\|_2 \leq \|A(x^* - x^\dagger)\|_2 + \|w\|_2,$$

which is at least $O(\sigma\sqrt{n})$ due to the term $\|w\|_2$. This upper-bound suggests that $\frac{1}{n} \|A(\hat{x} - x^\dagger)\|_2^2$ is bounded by $O(\epsilon^2\sigma^2) = O(\frac{\sigma^2 \mathcal{W}_1(\mathcal{K})^2}{m})$. This can be considered as a negative result for the sketching method, since the error scales as $O(\frac{1}{m})$ instead of $O(\frac{1}{n})$. We will show that this upper-bound is tight, and the $O(\frac{1}{m})$ scaling is unavoidable for all methods that sketch the data once. In contrast, as we will discuss in Section 1.5.7, an iterative sketching method can achieve optimal prediction error using sketches of comparable dimension.

We will in fact show that unless $m \geq n$, *any method* based on observing *only* the pair (SA, Sb) necessarily has a substantially larger error than the least-squares estimate. In particular, our result applies to an arbitrary measurable function $(SA, Sb) \mapsto \hat{x}$, which we refer to as an *estimator*.

More precisely, our lower bound applies to any random matrix $S \in \mathbb{R}^{m \times n}$ for which

$$\|\mathbb{E}[S^T(SS^T)^{-1}S]\|_{\text{op}} \leq \eta \frac{m}{n}, \quad (1.15)$$

where η is a constant independent of n and m , and $\|A\|_{\text{op}}$ denotes the ℓ_2 -operator norm, which reduces to the maximum eigenvalue for a symmetric matrix. These conditions hold for various standard choices of the sketching matrix, including most of those discussed in the Section 1.2: Gaussian sketch, ROS sketch⁶,

⁶ See [23] for a proof of this fact for Gaussian and ROS sketches. To be more precise, for ROS sketches, the condition (1.15) holds when rows are sampled without replacement.

the sparse sketch and α -balanced leverage sampling sketch. The following lemma shows that the condition (1.15) is satisfied for Gaussian sketches with equality, and $\eta = 1$.

LEMMA 1.4 *Let $S \in \mathbb{R}^{m \times n}$ be a random matrix with i.i.d. Gaussian entries. We have*

$$\|\mathbb{E}[S^T(SS^T)^{-1}S]\|_{op} = \frac{m}{n}.$$

Proof Let $S = U\Sigma V^T$ denote the Singular value decomposition of the random matrix S . Note that we have $S^T(SS^T)^{-1}S = VV^T$. By the rotation invariance of the Gaussian distribution, columns of V denoted by $\{v_i\}_{i=1}^m$ are uniformly distributed over the n dimensional unit sphere, and it holds that $\mathbb{E}v_iv_i^T = \frac{1}{n}I_n$ for $i = 1, \dots, m$. Consequently, we obtain

$$\mathbb{E}[S^T(SS^T)^{-1}S] = \mathbb{E}\sum_{i=1}^m v_iv_i^T = m\mathbb{E}v_1v_1^T = \frac{m}{n}I_n,$$

and the bound on the operator norm follows. \square

1.5.2 Fano's inequality

Let X and Y represent two random variables with a joint probability distribution $P_{x,y}$, where X is discrete and takes values from a finite set \mathcal{X} . Let $\hat{X} = g(Y)$ be the predicted value of X for some deterministic function g which also takes values in \mathcal{X} . Then Fano's inequality states that

$$P[X \neq \hat{X}] \geq \frac{H(X|Y) - 1}{\log_2(|\mathcal{X}| - 1)}.$$

Fano's inequality follows as a simple consequence of the chain rule for entropy. However it is very powerful for deriving lower bounds on the error probabilities in coding theory, statistics and machine learning [26, 27, 28, 29, 30, 7].

1.5.3 Metric entropy

For a given positive tolerance value $\delta > 0$, we define the δ -packing number $M_{\delta, \|\cdot\|}$ of a set $\mathcal{C} \subseteq \mathbb{R}^d$ with respect to a norm $\|\cdot\|$ as the largest number of vectors $\{x^j\}_{j=1}^M \subseteq \mathcal{C}$ which are elements of \mathcal{C} , and satisfy

$$\|x^k - x^l\| > \delta \quad \forall k \neq l.$$

We define the *metric entropy* of the set \mathcal{C} with respect to a norm $\|\cdot\|$, as the logarithm of the corresponding packing number

$$N_{\delta, \|\cdot\|}(\mathcal{C}) = \log_2 M_{\delta, \|\cdot\|}.$$

The concept of metric entropy provides a way to measure the complexity, or effective size of a set with infinitely many elements and dates back to the seminal work of Kolmogorov, Tikhomirov [31].

1.5.4 Minimax risk

In this chapter, we will take a frequentist approach in modeling the unknown vector x^\dagger we are trying to estimate from the data. In order to assess the quality of estimation, we will consider a risk function associated with our estimation method. Note that, for a fixed value of the unknown vector x^\dagger , there exists estimators which make no error for that particular vector x^\dagger , such as the estimator which always returns x^\dagger regardless of data. We will take the worst-case risk approach considered in the statistical estimation literature, which focus on the *minimax* risk. More precisely, we define the minimax risk as follows

$$\mathcal{M}(\mathcal{Q}) = \inf_{\hat{x} \in \mathcal{Q}} \sup_{x^\dagger \in \mathcal{X}} \mathbb{E} \left[\frac{1}{n} \|A(\hat{x} - x^\dagger)\|_2^2 \right], \quad (1.16)$$

where the infimum ranges over all estimators that use the input data A and b to estimate x^\dagger .

1.5.5 Reduction to hypothesis testing

In this section we present a reduction of minimax estimation risk to hypothesis testing. Suppose that we have a packing of the constraint set \mathcal{C} given by the collection $z^{(1)}, \dots, z^{(M)}$ with radius 2δ . More precisely, we have

$$\|A(z^{(i)} - z^{(j)})\|_2 \geq 2\delta \quad \forall i \neq j,$$

where $z^{(i)} \in \mathcal{C}$ for all $i = 1, \dots, M$. Next, consider a set of probability distributions $\{P_{z^{(j)}}\}_{j=1}^M$ corresponding to the distribution of the observation when the unknown vector is $x^\dagger = z^{(j)}$. Suppose that we have an M -ary hypothesis testing problem constructed as follows: Let J_δ denote a random variable with uniform distribution over the index set $\{1, \dots, M\}$ that allows us to pick an element of the packing set at random. Note that M is a function of δ , hence we keep the dependence of J_δ to δ explicit in our notation. Let us set the random variable Z according to the probability distribution $P_{z^{(j)}}$ in the event that $J_\delta = j$, i.e.,

$$Z \sim P_{z^{(j)}} \quad \text{whenever} \quad J_\delta = j.$$

Now we will consider the problem of detecting the index set given the value of Z . The next lemma is a standard reduction in minimax theory, and relates the minimax estimation risk to the M -ary hypothesis testing error (see Birge [30] and Yu [32]).

LEMMA 1.5 *The minimax risk \mathcal{Q} is lower bounded by*

$$\mathcal{M}(\mathcal{Q}) \geq \delta^2 \inf_{\psi} \mathbb{P}[\psi(Z) \neq J_\delta]. \quad (1.17)$$

A proof of this lemma can be found in Section 1.11. Lemma 1.5 allows us to apply Fano's method after transforming the estimation problem to an hypothesis

testing problem based on sketched data. Let us recall the condition on sketching matrices stated earlier

$$\|\mathbb{E}[S^T(SS^T)^{-1}S]\|_{\text{op}} \leq \eta \frac{m}{n}, \quad (1.18)$$

where η is a constant independent of n and m . Now we are ready to present the lower-bound on the statistical performance of sketching.

THEOREM 1.6 *For any random sketching matrix $S \in \mathbb{R}^{m \times n}$ satisfying condition (1.18), any estimator $(SA, Sb) \mapsto x^\dagger$ has MSE lower bounded as*

$$\sup_{x^\dagger \in \mathcal{C}_0} \mathbb{E}_{S,w} \left[\frac{1}{n} \|A(x^\dagger - x^*)\|_2^2 \right] \geq \frac{\sigma^2}{128\eta} \frac{\log_2(\frac{1}{2}M_{1/2})}{\min\{m, n\}} \quad (1.19)$$

where $M_{1/2}$ is the 1/2-packing number of $\mathcal{C}_0 \cap \mathbb{B}_A(1)$ in the semi-norm $\frac{1}{\sqrt{n}}\|A(\cdot)\|_2$.

We defer the proof to Section 1.8, and investigate the implications of the lower-bound in the next section. It can be shown that Theorem 1.6 is tight, since Theorem 1.1 provides a matching upper-bound.

1.5.6 Implications of the information-theoretic lower bound

We now investigate some consequences of the lower-bound given in Theorem 1.6. We will focus on concrete examples of popular statistical estimation and optimization problems to illustrate its applicability.

Example 1: Unconstrained least-squares

We first consider the simple unconstrained case, where the constraint is the entire d -dimensional space, i.e., $\mathcal{C} = \mathbb{R}^d$. With this choice, and for any data matrix A with, it is a classical result that under the observation model (1.13), the least-squares solution x^* has prediction mean-squared error upper bounded as follows⁷

$$\mathbb{E} \left[\frac{1}{n} \|A(x^* - x^\dagger)\|_2^2 \right] \lesssim \frac{\sigma^2 \text{rank}(A)}{n} \quad (1.20a)$$

$$\leq \frac{\sigma^2 d}{n}, \quad (1.20b)$$

where the expectation is over the noise variable w in (1.13). On the other hand, with the choice $\mathcal{C}_0 = \mathbb{B}_2(1)$, it is well known that we can construct a 1/2-packing with $M = 2^d$ elements, so that Theorem 1.6 implies that any estimator x^\dagger based on (SA, Sb) has prediction MSE lower bounded as

$$\mathbb{E}_{S,w} \left[\frac{1}{n} \|A(\hat{x} - x^\dagger)\|_2^2 \right] \gtrsim \frac{\sigma^2 d}{\min\{m, n\}}. \quad (1.20c)$$

⁷ In fact, a close form solution exists for the prediction error, which is straightforward to obtain from the closed form solution of the least squares estimator. However, this simple form is sufficient to illustrate information-theoretic lower-bounds.

Consequently, the sketch dimension m must grow proportionally to n in order for the sketched solution to have a mean-squared error comparable to the original least-squares estimate. This may not be desirable for least-squares problems in which $n \gg d$, since it should be possible to sketch down to a dimension proportional to $\text{rank}(A)$ which is always upper bounded by d . Thus, Theorem 1.6 reveals a surprising gap between the classical least-squares sketch (1.9) and the accuracy of the original least-squares estimate. In the regime $n \gg m$, the prediction MSE of the sketched solution is $O(\sigma^2 \frac{d}{m})$ which is a factor of $\frac{n}{m}$ larger than the optimal prediction MSE in (1.20b). In Section 1.5.7, we will see that this gap can be removed by iterative sketching algorithms which don't obey the information-theoretic lower-bound (1.20c).

Example 2: ℓ_1 constrained least-squares

We can consider other forms of constrained least-squares estimates as well, such as those involving an ℓ_1 -norm constraint to encourage sparsity in the solution. We now consider the sparse variant of the linear regression problem, which involves the ℓ_0 -“ball”

$$\mathbb{B}_0(s) := \{x \in \mathbb{R}^d \mid \sum_{j=1}^d \mathbb{I}[x_j \neq 0] \leq s\},$$

corresponding to the set of all vectors with at most s non-zero entries. Fixing some radius $R \geq \sqrt{s}$, consider a vector $x^\dagger \in \mathcal{C}_0 := \mathbb{B}_0(s) \cap \{\|x\|_1 = R\}$, and suppose that we have noisy observations the form $b = Ax^\dagger + w$.

Given this set-up, one way in which to estimate x^\dagger is by computing the least-squares estimate x^* constrained to the ℓ_1 -ball $\mathcal{C} = \{x \in \mathbb{R}^n \mid \|x\|_1 \leq R\}$ ⁸. This estimator is a form of the Lasso [33, 2] which is extensively studied in the context of statistical estimation and signal reconstruction.

On the other hand, the $\frac{1}{2}$ -packing number M of the set \mathcal{C}_0 can be lower bounded as $\log_2 M \gtrsim s \log_2 \left(\frac{ed}{s}\right)$. We refer the reader to [34] for a proof. Consequently, in application to this particular problem, Theorem 1.6 implies that any estimator \hat{x} based on the pair (SA, Sb) has mean-squared error lower bounded as

$$\mathbb{E}_{w,S} \left[\frac{1}{n} \|A(\hat{x} - x^\dagger)\|_2^2 \right] \gtrsim \frac{\sigma^2 s \log_2 \left(\frac{ed}{s}\right)}{\min\{m, n\}}. \quad (1.21a)$$

Again, we see that the projection dimension m must be of the order of n in order to match the mean-squared error of the constrained least-squares estimate x^* up to constant factors.

⁸ This set-up is slightly unrealistic, since the estimator is assumed to know the radius $R = \|x^\dagger\|_1$. In practice, one solves the least-squares problem with a Lagrangian constraint, but the underlying arguments are basically the same.

Example 3: Low rank matrix estimation

In the problem of multivariate regression, the goal is to estimate a matrix $X^\dagger \in \mathbb{R}^{d_1 \times d_2}$ model based on observations of the form

$$Y = AX^\dagger + W, \quad (1.22)$$

where $Y \in \mathbb{R}^{n \times d_2}$ is a matrix of observed responses, $A \in \mathbb{R}^{n \times d_1}$ is a data matrix, and $W \in \mathbb{R}^{n \times d_2}$ is a matrix of noise variables. A typical interpretation of this model is a collection of d_2 regression problems where each one involves a d_1 -dimensional regression vector, namely a particular column of the matrix X^\dagger . In many applications including reduced rank regression, multi-task learning and recommender systems (e.g., [35, 36, 37, 38]), it is reasonable to model the matrix X^\dagger as being low-rank. Note that a rank constraint on matrix X be written as an ℓ_0 -“norm” sparsity constraint on its singular values: in particular, we have

$$\text{rank}(X) \leq r \quad \text{if and only if} \quad \sum_{j=1}^{\min\{d_1, d_2\}} \mathbb{I}[\gamma_j(X) > 0] \leq r,$$

where $\gamma_j(X)$ denotes the j^{th} singular value of X . This observation motivates a standard relaxation of the rank constraint using the nuclear norm $\|X\|_{\text{nuc}} := \sum_{j=1}^{\min\{d_1, d_2\}} \gamma_j(X)$.

Accordingly, let us consider the constrained least-squares problem

$$X^* = \arg \min_{X \in \mathbb{R}^{d_1 \times d_2}} \left\{ \frac{1}{2} \|Y - AX\|_{\text{fro}}^2 \right\} \quad \text{such that } \|X\|_{\text{nuc}} \leq R, \quad (1.23)$$

where $\|\cdot\|_{\text{fro}}$ denotes the Frobenius norm on matrices, or equivalently the Euclidean norm on its vectorized version. Let \mathcal{C}_0 denote the set of matrices with rank $r < \frac{1}{2} \min\{d_1, d_2\}$, and Frobenius norm at most one. In this case the constrained least-squares solution X^* satisfies the bound

$$\mathbb{E} \left[\frac{1}{n} \|A(X^* - X^\dagger)\|_2^2 \right] \lesssim \frac{\sigma^2 r (d_1 + d_2)}{n}. \quad (1.24a)$$

On the other hand, the $\frac{1}{2}$ -packing number of the set \mathcal{C}_0 is lower bounded as $\log_2 M \gtrsim r(d_1 + d_2)$, (see [37] for a proof) so that Theorem 1.6 implies that any estimator \hat{X} based on the pair (SA, SY) has MSE lower bounded as

$$\mathbb{E}_{w,S} \left[\frac{1}{n} \|A(\hat{X} - X^\dagger)\|_2^2 \right] \gtrsim \frac{\sigma^2 r (d_1 + d_2)}{\min\{m, n\}}. \quad (1.24b)$$

As with the previous examples, we see the sub-optimality of the sketched approach in the regime $m < n$.

1.5.7 Iterative sketching

It is possible to improve the basic sketching estimator using adaptive measurements. Consider the constrained least squares problem in 1.8

$$x^* = \arg \min_{x \in \mathcal{C}} \frac{1}{2} \|Ax - b\|_2^2 \quad (1.25)$$

$$= \arg \min_{x \in \mathcal{C}} \underbrace{\frac{1}{2} \|Ax\|_2^2 - b^T Ax + \frac{1}{2} \|b\|_2^2}_{f(x)}. \quad (1.26)$$

We may use an iterative method to obtain x^* which uses the gradient $\nabla f(x) = A^T(Ax - b)$ and Hessian $\nabla^2 f(x) = A^T A$ to minimize the second order Taylor expansion of $f(x)$ at a current iterate x_t using $\nabla f(x_t)$ and $\nabla^2 f(x_t)$ as follows

$$x_{t+1} = x_t + \arg \min_{x \in \mathcal{C}} \left\| \left[\nabla^2 f(x) \right]^{1/2} x \right\|_2^2 + x^T \nabla f(x_t). \quad (1.27)$$

$$= x_t + \arg \min_{x \in \mathcal{C}} \|Ax\|_2^2 - x^T A^T (b - Ax_t). \quad (1.28)$$

We apply a sketching matrix S to the data A on the formulation (1.28) and define this procedure as *Iterative Sketch*

$$x_{t+1} = x_t + \arg \min_{x \in \mathcal{C}} \|SAx\|_2^2 - 2x^T A^T (b - Ax_t). \quad (1.29)$$

Note that this procedure uses more information than the classical sketch 1.9, in particular it calculates the left matrix-vector multiplies with the data A in the following order:

$$\begin{aligned} & s_1^T A \\ & s_2^T A \\ & s_m^T A \\ & \vdots \\ & (b - Ax_1)^T A \\ & \vdots \\ & (b - Ax_t)^T A, \end{aligned}$$

where s_1^T, \dots, s_m^T are the rows of the sketching matrix S . This can be considered as an adaptive form of sketching where the residual directions $(b - Ax_t)$ are used after the random directions s_1, \dots, s_m . As a consequence, the information-theoretic bounds we considered in Section 1.4.6 do not apply to iterative sketching. In Pilanci and Wainwright [23], it is shown that this algorithm achieves the minimax statistical given in (1.16) using at most $O(\log_2 n)$ iterations while obtaining equivalent speedups from sketching. We also note that iterative sketching method can also be applied to more general convex optimization problems other than the least squares objective. We refer the reader to Pilanci and Wainwright [39] for the application of sketching in solving general convex optimization problems.

1.6 Nonparametric problems

1.6.1 Nonparametric regression

In this section we discuss an extension of the sketching method to nonparametric regression problems over Hilbert spaces. The goal of non-parametric regression is making predictions of a continuous response after observing a covariate where they are related via

$$y_i = f^*(x_i) + w_i, \quad \text{for} \quad (1.30)$$

where $w \sim \mathcal{N}(0, \sigma^2 I_n)$, and the function $f^*(x)$ needs to be estimated based on $\{x_i, y_i\}_{i=1}^n$. We will consider the well-studied case where the function f^* is assumed to belong a reproducing kernel Hilbert space (RKHS) \mathcal{H} , and has a bounded Hilbert norm $\|f\|_{\mathcal{H}}$ [40, 41]. For these regression problems it is customary to consider the kernel ridge regression (KRR) problem based on convex optimization

$$\hat{f} = \arg \min_{f \in \mathcal{H}} \left\{ \frac{1}{2n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{H}}^2 \right\} \quad (1.31)$$

An RKHS is generated by a kernel function which positive semidefinite (PSD). A PSD kernel is a symmetric function $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ that satisfies

$$\sum_{i,j=1}^r y_i y_j \mathcal{K}(x_i, x_j) \geq 0$$

for all collections of points $\{x_1, \dots, x_n\}$, $y \in \mathbb{R}^n$ and $\forall r \in \mathbb{Z}_+$. The vector space of all functions of the form

$$f(\cdot) = \sum_i^r y_i \mathcal{K}(\cdot, x_i),$$

generates an RKHS by taking closure of all such linear combinations. It can be shown that this RKHS is uniquely associated with the kernel function \mathcal{K} (see Aronszajn [42]). Let us define a finite dimensional kernel matrix using n covariates as follows

$$K_{ij} = \frac{1}{n} \mathcal{K}(x_i, x_j),$$

which is a positive semidefinite matrix. In the linear least squares regression the kernel matrix reduces to the Gram matrix given by $K = A^T A$. It is also known that the above infinite dimensional program can be recast as a finite dimensional quadratic optimization problem involving the kernel matrix

$$\hat{w} = \arg \min_{w \in \mathbb{R}^n} \frac{1}{2} \|Kw - \frac{1}{\sqrt{n}} y\|_2^2 + \lambda w^T K w \quad (1.32)$$

$$= \arg \min_{w \in \mathbb{R}^n} \frac{1}{2} w^T K^2 w - w^T \frac{Ky}{\sqrt{n}} + \lambda w^T K w, \quad (1.33)$$

and we can find the optimal solution to the infinite dimensional problem (1.31) via the following relation⁹

$$\hat{f}(\cdot) = \frac{1}{n} \sum_{i=1}^n \hat{w}_i \mathcal{K}(\cdot, x_i). \quad (1.34)$$

We now define a kernel complexity measure based on the eigenvalues of the kernel matrix K . Let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ correspond to the real eigenvalues of the symmetric positive definite kernel matrix K . The kernel complexity is defined as follows.

DEFINITION 1.7 (Kernel complexity)

$$\mathcal{R}(\delta) = \sqrt{\sum_{i=1}^n \min\{\delta^2, \lambda_i\}}$$

which is the sum of eigenvalues truncated at level δ . Similar to (1.14), we define a critical radius $\delta^*(n)$ equal the smallest positive solution $\delta^*(n) > 0$ to the following inequality

$$\frac{\mathcal{R}(\delta)}{\delta\sqrt{n}} \leq \frac{\delta}{\sigma}, \quad (1.35)$$

where σ is the noise standard deviation in the statistical model (1.30). The existence of a unique is guaranteed for all kernel classes (see Bartlett et al. [20]). The critical radius plays an important role in the minimax risk through an information-theoretic argument. The next theorem provides a lower-bound on the statistical risk of any estimator applied to the observation model (1.30).

THEOREM 1.8 *Given n i.i.d. samples from the model (1.30), any estimator \hat{f} has prediction error lower bounded as*

$$\sup_{\|f^*\|_{\mathcal{H}} \leq 1} \mathbb{E} \frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f^*(x_i))^2 \geq c_0 \delta^*(n)^2, \quad (1.36)$$

where c_0 is a numerical constant and $\delta^*(n)$ is the critical radius defined in (1.35)

The lower-bound given by Theorem 1.8 can be shown to be tight, and achieved by the kernel based optimization procedure (1.33) and (1.34) (see Bartlett et al. [20]). The proof of Theorem 1.8 can be found in Yang et al. [43]. We may define the *effective dimension* $d^*(n)$ of the kernel via the relation

$$d^*(n) := n\delta^*(n)^2,$$

This definition allows us to interpret the convergence rate in (1.36) as $\frac{d^*(n)}{n}$,

⁹ Our definition of the kernel optimization problem slightly differs from the literature. The classical kernel problem can be recovered by a variable change $w' = K^{1/2}w$ where $K^{1/2}$ is the matrix square-root. We refer the reader to [41] for more details on kernel based methods.

which resembles the classical parametric convergence rate where the number of variables is $d^*(n)$.

1.6.2 Sketching kernels

Solving the optimization problem (1.33) becomes a computational challenge when the sample size n is large since it involves linear algebraic operations on an $n \times n$ matrix K . There is a large body of literature on approximating kernel matrices using randomized methods [44, 45, 46, 47]. Here we assume that the matrix K is available, and a sketching matrix $S \in \mathbb{R}^{m \times n}$ can be applied to form a randomized approximation of the kernel matrix. We will present an extension of (1.9), which achieves optimal statistical accuracy. Specifically, the sketching method we consider solves

$$\hat{v} = \arg \min_{v \in \mathbb{R}^m} \frac{1}{2} v^T (SK)(KS^T)v - v^T \frac{SKy}{\sqrt{n}} + \lambda v^T SKS^T v, \quad (1.37)$$

which involves smaller dimensional sketched kernel matrices SK , SKS^T and a lower dimensional decision variable $v \in \mathbb{R}^m$. Then we can recover the original variable via $w = S^T v$. The next theorem shows that the sketched kernel based optimization method achieves optimal prediction error.

THEOREM 1.9 *Let $S \in \mathbb{R}^{m \times n}$ be a Gaussian sketching matrix where $m \geq c_3 d_n$, and choose $\lambda = 3\delta^*(n)$. Given n i.i.d. samples from the model (1.30), the sketching procedure (1.42) produces a regression estimate \hat{f} which satisfies the bound*

$$\frac{1}{n} \sum_{i=1}^n (\hat{f}(x_i) - f^*(x_i))^2 \leq c_2 \delta^*(n)^2,$$

where $\delta^*(n)$ is the critical radius defined in (1.35)

A proof of this theorem can be found in Yang et al. [43]. We note that a similar result holds for the ROS sketch matrices with extra logarithmic terms in the dimension of the sketch, i.e., when $m \geq c_4 d_n \log^4(n)$ holds. Notably, Theorem 1.9 guarantees that the sketched estimator achieves the optimal error. This is in contrast to the lower-bound Section 1.4.6 where the sketching method does not achieve a minimax optimal error. This is due to the fact that the sketched problem in (1.37) is using the observation SKy instead of Sy . Therefore, the lower bound in Section 1.4.6 does not apply for this construction. It is worth noting that one can formulate the ordinary least squares as a kernel regression problem with kernel $K = AA^T$, and then apply the sketching method (1.37), which is guaranteed to achieve minimax optimal risk. However, computing the kernel matrix AA^T would cost $O(nd^2)$ operations, which is larger than solving the original least squares problem.

We note that some kernel approximation methods avoid computing the kernel matrix K and directly forms low rank approximations. We refer the reader to

[44] for an example, which also provides an error guarantee for the approximate kernel.

1.7 Extensions: Privacy and communication complexity

1.7.1 Privacy and information-theoretic bounds

Another interesting property of randomized sketching is privacy preservation in the context of optimization and learning. Privacy properties of random projections for various statistical tasks are studied in the recent literature [9],[48], [10]. It is of great theoretical and practical interest to characterize fundamental privacy and optimization trade-offs of randomized algorithms. We first show the relation between sketching and a mutual information based privacy measure. Suppose that we

1.7.2 Mutual information privacy

Suppose we model the data matrix $A \in \mathbb{R}^{n \times d}$ as stochastic where each entry is drawn randomly. One way we can assess the revealed information to the server is considering mutual information per symbol, which is given by the formula

$$\begin{aligned} \frac{I(SA; A)}{nd} &= \frac{1}{nd} \{H(A) - H(A|SA)\} \\ &= \frac{1}{nd} D(\mathbb{P}_{SA,A} || \mathbb{P}_{SA,A}) \end{aligned}$$

where we normalize by nd since the data matrix A has nd entries in total. The following corollary is a direct application of Theorem 1.1.

COROLLARY *Let the entries of the matrix A be i.i.d from an arbitrary distribution with finite variance σ^2 . Using sketched data, we can obtain an ϵ -approximate¹⁰ solution to the optimization problem while ensuring that revealed the mutual information satisfies*

$$\frac{I(SA; A)}{nd} \leq \frac{c_0}{\epsilon^2} \frac{\mathcal{W}^2(AK)}{n} \log_2(2\pi e \sigma^2).$$

Therefore, we can guarantee mutual information privacy of the sketching based methods, whenever the term $\mathcal{W}(AK)$ is small.

An alternative and popular characterization of privacy is referred as the differential privacy (see Dwork et al, [9]), where other randomized methods such as additive noise for preserving privacy was studied. It is also possible to directly analyze differential privacy preserving aspects of random projections as considered in Blocki et al. [10].

¹⁰ ϵ -approximate solution refers to the approximation defined in Theorem 1.1, relative to the optimal value.

1.7.3 Optimization based privacy attacks

We briefly discuss a possible approach an adversary might take to circumvent the privacy provided by sketching. If the data matrix is sparse, then one might consider optimization based recovery techniques borrowed from Compressed Sensing to recover the data A given the sketched data $\tilde{A} = SA$.

$$\begin{aligned} \min_A \|A\|_1 \\ \text{s.t. } SA = \tilde{A} \end{aligned}$$

The success of the above optimization method will critically depend on the sparsity level of the original data A . Most of the randomized sketching constructions shown in Section 1.2 can be shown to be susceptible to data recovery via optimization (see Candès and Tao [25], and Candès et al. [49]). However, this method assumes that the sketching matrix S is available to the attacker. If the S is not available to the adversary, then the above method cannot be used and the recovery is not straightforward.

1.7.4 Communication complexity space lower bounds

In this section we consider a streaming model of computation, where the algorithm is allowed to make only one pass over the data. In this model, an algorithm receives updates to the entries of the data matrix A in the form “add a to A_{ij} ”. An entry can be updated more than once, and the value a is any arbitrary real number. The sketches introduced in this chapter provide a valuable data structure when the matrix is very large in size, and storing and updating the matrix directly can be impractical. Due to the linearity of sketches, we can update the sketch SA by adding $Sae_i e_j^T$ to SA , and maintain an approximation with limited memory.

The following theorem due to Clarkson and Woodruff [50], provides a lower-bound of the space used by any algorithm for least squares regression which performs a single pass over the data.

THEOREM 1.10 *Any randomized 1-pass algorithm which returns an ϵ -approximate solution to the unconstrained least squares problem with probability at least $7/9$ needs $\Omega(d^2(\frac{1}{\epsilon} + \log(nd)))$ bits of space.*

This theorem confirms that the space complexity of sketching for unconstrained least squares regression is near optimal. By the choice of the sketching dimension $m = O(d)$, the space used by the sketch SA is $O(d^2)$ which is optimal up to constants according to the theorem.

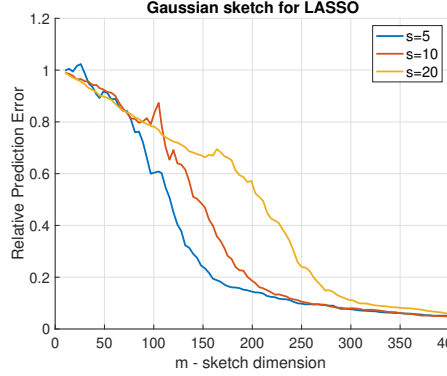


Figure 1.5 Sketching LASSO using Gaussian random projections

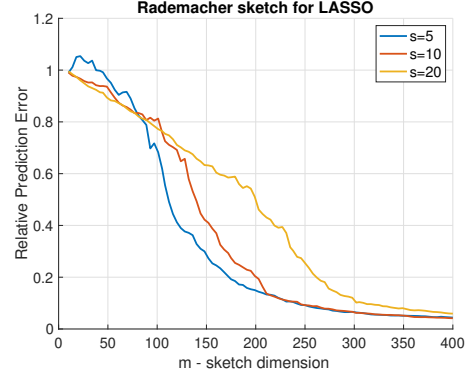


Figure 1.6 Sketching LASSO using Rademacher random projections

1.8 Numerical experiments

In this section we illustrate the sketching method numerically and confirm the theoretical predictions of Theorems 1.1 and 1.6. We consider both the classical low-dimensional statistical regime where $n > d$, and the ℓ_1 constrained least squares minimization known as LASSO (see Tibshirani [51]):

$$x^* = \arg \min_{x \text{ s.t. } \|x\|_1 \leq \lambda} \|Ax - b\|_2.$$

We generate a random i.i.d. data matrix $A \in \mathbb{R}^{n \times d}$ where $n = 10000$ and $d = 1000$, and set the observation vector $b = Ax^\dagger + \sigma w$, where $x^\dagger \in \{-1, 0, 1\}^d$ is a random s -sparse vector, w has i.i.d. $\mathcal{N}(0, 10^{-4})$ components. For the sketching matrix $S \in \mathbb{R}^{m \times n}$, we consider the Gaussian and Rademacher (± 1 i.i.d. valued) random matrices where m ranges between 10 and 400. Consequently, we solve the sketched program

$$\hat{x} = \arg \min_{x \text{ s.t. } \|x\|_1 \leq \lambda} \|SAx - Sb\|_2.$$

Figure 1.5 and 1.6 shows the relative prediction mean squared error given by the ratio

$$\frac{\frac{1}{n} \|A(\hat{x} - x^\dagger)\|_2^2}{\frac{1}{n} \|A(x^* - x^\dagger)\|_2^2},$$

where it is averaged over 20 realizations of the sketching matrix, \hat{x} and x^* are the sketched and the original solutions respectively. As predicted by the upper and lower bounds given in Theorem 1.1 and Theorem 1.6, the prediction mean squared error of the sketched estimator scales as $O\left(\frac{s \log d}{m}\right)$, since the corresponding Gaussian complexity $\mathcal{W}_1(\mathcal{K})^2$ is $O(s \log d)$. These plots reveal that the prediction mean squared error of the sketched estimators for both Gaussian and Rademacher sketches are in agreement with the theory.

1.9 Summary

This chapter presented an overview of random projection based methods for solving large scale statistical estimation and constrained optimization problems. We investigated fundamental lower-bounds on the performance of sketching using information theoretical tools. Randomized sketching has interesting theoretical properties, and also has numerous practical advantages in machine learning and optimization problems. Sketching yields faster algorithms with lower space complexity while maintaining strong approximation guarantees.

For the upper-bound on the approximation accuracy in Theorem 1.3, Gaussian complexity play an important role, and also provide a geometric characterization of the dimension of the sketch. The lower-bounds given in Theorem 1.6 is statistical in nature, and involve packing numbers, and consequently metric entropy which measure the complexity of the sets. The upper bounds on the Gaussian sketch can be extended to Rademacher sketches, sub-Gaussian sketches, and Randomized Orthogonal System sketches (see Pilanci and Wainwright, and also Yun et al. [22, 43] for the proofs). However, the results for non-Gaussian sketches often involve superfluous logarithmic factors and large constants as artifacts of the analysis. As it can be observed in Figure 1.5 and Figure 1.6, the mean squared error curves for Gaussian and Rademacher sketches are in agreement with each other. It can be conjectured that the approximation ratio of sketching is universal for random matrices with entries sampled from well behaved distributions. This is an important theoretical question for future research. We refer the reader to the work of Donoho and Tanner [52] for the observations of universality in compressed sensing.

Finally, a number of important limitations of the analysis techniques need to be considered. The minimax criteria (1.16) is worst case in nature by its definition, and may not correctly reflect the average error of sketching when the unknown vector x^\dagger is randomly distributed. Furthermore, in some applications, it might be suitable to consider prior information on the unknown vector. As an interesting direction of future research, it would be interesting to study lower bounds for sketching in a Bayesian setting.

1.10 Proof of Theorem 1.6

Let us define the short-hand notation $\|\cdot\|_A := \frac{1}{\sqrt{n}}\|A(\cdot)\|_2$. Let $\{z^j\}_{j=1}^M$ be a $1/2$ -packing of $\mathcal{C}_0 \cap \mathbb{B}_A(1)$ in the semi-norm defined by $\|\cdot\|_A$, and for a fixed $\delta \in (0, 1/4)$, define $x^j = 4\delta z^j$. Since $4\delta \in (0, 1)$, the star-shaped assumption guarantees that each x^j belongs to \mathcal{C}_0 . We thus obtain a collection of M vectors in \mathcal{C}_0 such that

$$2\delta \leq \|x^j - x^k\|_A \leq 8\delta \quad \text{for all } j \neq k.$$

Letting J be a random index uniformly distributed over $\{1, \dots, M\}$, suppose that conditionally on $J = j$, we observe the sketched observation vector $Sb = SAx^j + Sw$, as well as the sketched matrix SA . Conditioned on $J = j$, the random vector Sb follows a $\mathcal{N}(SAx^j, \sigma^2 SS^T)$ distribution, denoted by \mathbb{P}_{x^j} . We let \bar{Y} denote the resulting mixture variable, with distribution $\frac{1}{M} \sum_{j=1}^M \mathbb{P}_{x^j}$.

Consider the multi-way testing problem of determining the index J based on observing \bar{Y} . With this set-up, we may apply Lemma 1.5 (see e.g., [30, 32]) which implies that, for any estimator x^\dagger , the worst-case mean-squared error is lower bounded as

$$\sup_{x^* \in \mathcal{C}} \mathbb{E}_{S,w} \|x^\dagger - x^*\|_A^2 \geq \delta^2 \inf_{\psi} \mathbb{P}[\psi(\bar{Y}) \neq J], \quad (1.38)$$

where the infimum ranges over all testing functions ψ . Consequently, it suffices to show that the testing error is lower bounded by $1/2$.

In order to do so, we first apply Fano's inequality [27] conditionally on the sketching matrix S and get

$$\mathbb{P}[\psi(\bar{Y}) \neq J] = \mathbb{E}_S \left\{ \mathbb{P}[\psi(\bar{Y}) \neq J \mid S] \right\} \geq 1 - \frac{\mathbb{E}_S [I_S(\bar{Y}; J)] + 1}{\log_2 M}, \quad (1.39)$$

where $I_S(\bar{Y}; J)$ denotes the mutual information between \bar{Y} and J with S fixed. Our next step is to upper bound the expectation $\mathbb{E}_S [I(\bar{Y}; J)]$.

Letting $D(\mathbb{P}_{x^j} \parallel \mathbb{P}_{x^k})$ denote the Kullback-Leibler divergence between the distributions \mathbb{P}_{x^j} and \mathbb{P}_{x^k} , the convexity of Kullback-Leibler divergence implies that

$$\begin{aligned} I_S(\bar{Y}; J) &= \frac{1}{M} \sum_{j=1}^M D(\mathbb{P}_{x^j} \parallel \frac{1}{M} \sum_{k=1}^M \mathbb{P}_{x^k}) \\ &\leq \frac{1}{M^2} \sum_{j,k=1}^M D(\mathbb{P}_{x^j} \parallel \mathbb{P}_{x^k}). \end{aligned}$$

Computing the KL divergence for Gaussian vectors yields

$$I_S(\bar{Y}; J) \leq \frac{1}{M^2} \sum_{j,k=1}^M \frac{1}{2\sigma^2} (x^j - x^k)^T A^T \left[S^T (SS^T)^{-1} S \right] A (x^j - x^k).$$

Thus, using condition (1.15), we have

$$\mathbb{E}_S[I(\bar{Y}; J)] \leq \frac{1}{M^2} \sum_{j,k=1}^M \frac{m \eta}{2 n \sigma^2} \|A(x^j - x^k)\|_2^2 \leq \frac{32 m \eta}{\sigma^2} \delta^2,$$

where the final inequality uses the fact that $\|x^j - x^k\|_A = \frac{1}{\sqrt{n}} \|A(x^j - x^k)\|_2 \leq 8\delta$ for all pairs.

Combined with our previous bounds (1.38) and (1.39), we find that

$$\sup_{x^* \in \mathcal{C}} \mathbb{E} \|\hat{x} - x^*\|_2^2 \geq \delta^2 \left\{ 1 - \frac{32 m \eta \delta^2}{\sigma^2} + 1 \right\}.$$

Setting $\delta = \frac{\sigma^2 \log_2(M/2)}{64 \eta m}$ yields the lower bound (1.19).

1.11 Proof of Lemma 1.5

Proof By Markov's inequality applied on the random variable $\|\hat{x} - x^\dagger\|_A^2$ we have

$$\mathbb{E} \|\hat{x} - x^\dagger\|_A^2 \geq \delta^2 \mathbb{P}[\|\hat{x} - x^\dagger\|_A^2 \geq \delta^2] \quad (1.40)$$

Now note that

$$\begin{aligned} \sup_{x^* \in \mathcal{C}} \mathbb{P}[\|\hat{x} - x^\dagger\|_A \geq \delta] &\geq \max_{j \in \{1, \dots, M\}} \mathbb{P}[\|\hat{x} - x^{(j)}\|_A \geq \delta \mid J_\delta = j] \\ &\geq \frac{1}{M} \sum_{j=1}^M \mathbb{P}[\|\hat{x} - x^{(j)}\|_A \geq \delta \mid J_\delta = j], \end{aligned} \quad (1.41)$$

since every element of the packings set satisfies $x^{(j)} \in \mathcal{C}$ and the discrete maximum is upper bounded by the average over $\{1, \dots, M\}$. Since we have $\mathbb{P}[J_\delta = j] = \frac{1}{M}$, we equivalently have

$$\begin{aligned} \frac{1}{M} \sum_{j=1}^M \mathbb{P}[\|\hat{x} - x^{(j)}\|_A \geq \delta \mid J_\delta = j] &= \sum_{j=1}^M \mathbb{P}[\|\hat{x} - x^{(j)}\|_A \geq \delta \mid J_\delta = j] \mathbb{P}[J_\delta = j] \\ &= \mathbb{P}[\|\hat{x} - x^{(J_\delta)}\|_A \geq \delta] \end{aligned} \quad (1.42)$$

Now we will argue that, whenever the true index is $J_\delta = j$ and if $\|\hat{x} - x^{(j)}\|_A < \delta$, then we can form a hypothesis tester $\psi(Z)$ identifying the true index j . Consider the test

$$\psi(Z) := \arg \min_{j \in [M]} \|x^{(j)} - \hat{x}\|_A.$$

Now note that $\|x^j - \hat{x}\|_A < \delta$ ensures that

$$\|x^{(i)} - \hat{x}\|_A \geq \|x^{(i)} - x^{(j)}\|_A - \|x^{(j)} - \hat{x}\|_A \geq 2\delta - \delta = \delta,$$

where the second inequality follows from the 2δ -packing construction of our collection $x^{(1)}, \dots, x^{(M)}$. Consequently $\|x^{(i)} - \hat{x}\|_A > \delta$ for all $i \in \{1, \dots, M\} - \{j\}$, and the test $\psi(Z)$ identifies the true index $J = j$. Therefore we obtain

$$\left\{ \|x^{(j)} - \hat{x}\|_A < \delta \right\} \implies \{ \phi(Z) = j \},$$

and conclude that the complements of these events obey

$$\mathbb{P} \left[\|x^{(j)} - \hat{x}\|_A \geq \delta \mid J_\delta = j \right] \geq \mathbb{P} [\phi(Z) \neq j \mid J_\delta = j].$$

Taking averages over the indices $1, \dots, M$ we obtain

$$\mathbb{P} \left[\|x^{(J_\delta)} - \hat{x}\|_A \geq \delta \right] = \frac{1}{M} \sum_{j=1}^M \mathbb{P} \left[\|x^{(j)} - \hat{x}\|_A \geq \delta \mid J_\delta = j \right] \geq \mathbb{P} [\phi(Z) \neq J_\delta].$$

Combining the above with the earlier lower-bound (1.41) and the identity (1.42), we obtain

$$\sup_{x^* \in \mathcal{C}} \mathbb{P} [\|\hat{x} - x^*\|_A \geq \delta] \geq \mathbb{P} [\phi(Z) \neq J_\delta] \geq \inf_{\phi} \mathbb{P} [\phi(Z) \neq J_\delta]$$

where the second inequality follows by taking the infimum over all tests, which can only make the probability smaller. Plugging in the above lower-bound in (1.40) completes the proof of the lemma. \square

References

- [1] S. Vempala. *The Random Projection Method*. Discrete Mathematics and Theoretical Computer Science. American Mathematical Society, Providence, RI, 2004.
- [2] E. J. Candes and T. Tao. Near-optimal signal recovery from random projections: Universal encoding strategies? *IEEE Trans. Information Theory*, 52(12):5406–5425, 2006.
- [3] N. Halko, P. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 53(2):217–288, 2011.
- [4] M. W. Mahoney. Randomized algorithms for matrices and data. *Foundations and Trends in Machine Learning in Machine Learning*, 3(2), 2011.
- [5] D. P. Woodruff. Sketching as a tool for numerical linear algebra. *Foundations and Trends® in Theoretical Computer Science*, 10(1–2):1–157, 2014.
- [6] S. Muthukrishnan. Data streams: Algorithms and applications. *Foundations and Trends® in Theoretical Computer Science*, 1(2):117–236, 2005.
- [7] B. Yu. Assouad, Fano and Le Cam. In *Festschrift in Honor of L. Le Cam on his 70th Birthday*. 1993.
- [8] Y. Yang and A. Barron. Information-theoretic determination of minimax rates of convergence. *Annals of Statistics*, 27(5):1564–1599, 1999.
- [9] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography Conference*, pages 265–284. Springer, 2006.
- [10] J. Blocki, A. Blum, A. Datta, and O. Sheffet. The Johnson-Lindenstrauss transform itself preserves differential privacy. In *2012 IEEE 53rd Annual Symposium on Foundations of Computer Science*, pages 410–419, Oct 2012.
- [11] N. Ailon and B. Chazelle. Approximate nearest neighbors and the fast Johnson-Lindenstrauss transform. In *Proceedings of the 38th Annual ACM Symposium on Theory of Computing*, pages 557–563. ACM, 2006.
- [12] P. Drineas and M. W. Mahoney. Effective resistances, statistical leverage, and applications to linear equation solving. *arXiv preprint arXiv:1005.3097*, 2010.
- [13] D. A. Spielman and N. Srivastava. Graph sparsification by effective resistances. *SIAM Journal on Computing*, 40(6):1913–1926, 2011.
- [14] Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams. In *International Colloquium on Automata, Languages, and Programming*, pages 693–703. Springer, 2002.
- [15] D. M. Kane and J. Nelson. Sparser Johnson-Lindenstrauss transforms. *Journal of the ACM*, 61(1), 2014.

-
- [16] Jelani Nelson and Huy L. Nguyễn. 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.
- [17] J. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms*, volume 1. Springer-Verlag, New York, 1993.
- [18] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, Cambridge, UK, 2004.
- [19] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*. Springer-Verlag, New York, NY, 1991.
- [20] P. L. Bartlett, O. Bousquet, and S. Mendelson. Local Rademacher complexities. *Annals of Statistics*, 33(4):1497–1537, 2005.
- [21] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky. The convex geometry of linear inverse problems. *Foundations of Computational Mathematics*, 12(6):805–849, 2012.
- [22] M. Pilanci and M. J. Wainwright. Randomized sketches of convex programs with sharp guarantees. Technical report, UC Berkeley, 2014. Full length version at arXiv:1404.7203; Presented in part at ISIT 2014.
- [23] M. Pilanci and M. J. Wainwright. Iterative hessian sketch: Fast and accurate solution approximation for constrained least-squares. *The Journal of Machine Learning Research*, 17(1):1842–1879, 2016.
- [24] S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM J. Sci. Computing*, 20(1):33–61, 1998.
- [25] E. J. Candes and T. Tao. Decoding by linear programming. *IEEE Trans. Information Theory*, 51(12):4203–4215, December 2005.
- [26] Robert M Fano and WT Wintringham. Transmission of information. *Physics Today*, 14:56, 1961.
- [27] T.M. Cover and J.A. Thomas. *Elements of Information Theory*. John Wiley and Sons, New York, 1991.
- [28] P. Assouad. Deux remarques sur l’estimation. *Comptes Rendus de l’Academie des Sciences de Paris*, 296:1021–1024, 1983.
- [29] I. A. Ibragimov and R. Z. Has’minskii. *Statistical Estimation: Asymptotic Theory*. Springer-Verlag, New York, 1981.
- [30] L. Birgé. Estimating a density under order restrictions: Non-asymptotic minimax risk. *Annals of Statistics*, 15(3):995–1012, March 1987.
- [31] A. Kolmogorov and B. Tikhomirov. ϵ -entropy and ϵ -capacity of sets in functional spaces. *Uspekhi Mat. Nauk.*, 86:3–86, 1959. Appeared in English as Amer. Math. Soc. Translations, 17:277–364, 1961.
- [32] B. Yu. Assouad, Fano and Le Cam. In *Festschrift for Lucien Le Cam*, pages 423–435. Springer-Verlag, Berlin, 1997.
- [33] R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*, 58(1):267–288, 1996.
- [34] G. Raskutti, M. J. Wainwright, and B. Yu. Minimax rates of estimation for high-dimensional linear regression over ℓ_q -balls. *IEEE Trans. Information Theory*, 57(10):6976–6994, October 2011.
- [35] N. Srebro, N. Alon, and T. S. Jaakkola. Generalization error bounds for collaborative prediction with low-rank matrices. In *Advances In Neural Information Processing Systems*, pages 1321–1328, 2005.

-
- [36] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society B*, 1(68):49, 2006.
- [37] S. Negahban and M. J. Wainwright. Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *Annals of Statistics*, 39(2):1069–1097, 2011.
- [38] F. Bunea, Y. She, and M. Wegkamp. Optimal selection of reduced rank estimators of high-dimensional matrices. *Annals of Statistics*, 39(2):1282–1309, 2011.
- [39] M. Pilanci and M. J. Wainwright. Newton sketch: A near linear-time optimization algorithm with linear-quadratic convergence. *SIAM Journal on Optimization*, 27(1):205–245, 2017.
- [40] Howard L. Weinert (ed.), editor. *Reproducing Kernel Hilbert Spaces : Applications in Statistical Signal Processing*. Hutchinson Ross Publishing Co., Stroudsburg, PA, 1982.
- [41] B. Schölkopf and A. Smola. *Learning with Kernels*. MIT Press, Cambridge, MA, 2002.
- [42] N. Aronszajn. Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68:337–404, 1950.
- [43] Y. Yang, M. Pilanci, M. J. Wainwright, et al. Randomized sketches for kernels: Fast and optimal nonparametric regression. *The Annals of Statistics*, 45(3):991–1023, 2017.
- [44] A. Rahimi and B. Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2008.
- [45] A. Rahimi and B. Recht. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In *Advances in neural information processing systems*, pages 1313–1320, 2009.
- [46] P. Drineas and M. W. Mahoney. On the nyström method for approximating a gram matrix for improved kernel-based learning. *journal of machine learning research*, 6(Dec):2153–2175, 2005.
- [47] Q. Le, T. Sarlós, and A. Smola. Fastfood-approximating kernel expansions in loglinear time. In *Proceedings of the international conference on machine learning*, volume 85, 2013.
- [48] S. Zhou, J. Lafferty, and L. Wasserman. Compressed regression. In *Neural Information Processing Systems*, December 2007.
- [49] E. J. Candes, J. Romberg, and T. Tao. Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information. *IEEE Trans. Information Theory*, 52(2):489–509, February 2004.
- [50] K. L. Clarkson and D. 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.
- [51] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- [52] David Donoho and Jared Tanner. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 367(1906):4273–4293, 2009.