

Vol. 50, No. 4, November 2025, pp. 2941–2971 ISSN 0364-765X (print), ISSN 1526-5471 (online)

Quadratic Memory Is Necessary for Optimal Query Complexity in Convex Optimization: Center of Mass Is Pareto Optimal

Moïse Blanchard, a,* Junhui Zhang, Patrick Jailleta

^a Operations Research Center, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139

Received: July 8, 2023
Revised: July 11, 2024
Accepted: September 5, 2024
Rublished Online in Articles in Ad

Published Online in Articles in Advance: November 20, 2024

MSC2020 Subject Classifications: Primary: 68Q25

00020

https://doi.org/10.1287/moor.2023.0208

Copyright: © 2024 INFORMS

Abstract. We give query complexity lower bounds for convex optimization and the related feasibility problem. We show that quadratic memory is necessary to achieve the optimal oracle complexity for first-order convex optimization. In particular, this shows that center-of-mass cutting-plane algorithms in dimension d, which use $\tilde{O}(d^2)$ memory and $\tilde{O}(d)$ queries, are Pareto optimal for both convex optimization and the feasibility problem, up to logarithmic factors. Precisely, building upon techniques introduced in previous works, we prove that to minimize 1-Lipschitz convex functions over the unit ball to $1/d^4$ accuracy, any deterministic first-order algorithms using at most $d^{2-\delta}$ bits of memory must make $\tilde{\Omega}(d^{1+\delta/3})$ queries for any $\delta \in [0,1]$. For the feasibility problem, in which an algorithm only has access to a separation oracle, we show a stronger trade-off; for at most $d^{2-\delta}$ memory, the number of queries required is $\tilde{\Omega}(d^{1+\delta})$. This resolves a Conference on Learning Theory 2019 open problem.

Funding: This work was partly supported by the Air Force Office of Scientific Research [Grant FA9550-19-1-0263] and the Office of Naval Research [Grant N00014-18-1-2122].

Keywords: convex optimization • feasibility problem • first-order methods • cutting plane • center of mass • memory lower bounds • query complexity

1. Introduction

We consider the canonical problem of first-order convex optimization, in which one aims to minimize a convex function $f: \mathbb{R}^d \to \mathbb{R}$ with access to an oracle that for any query x, returns $(f(x), \nabla f(x))$ the value of the function and a subgradient of f at x. Arguably, this is one of the most fundamental problems in optimization, mathematical programming, and machine learning.

A classical question is how many oracle queries are required to guarantee finding an ϵ -approximate minimizer for any 1-Lipschitz convex functions $f: \mathbb{R}^d \to \mathbb{R}$ over the unit ball. We denote by $B_d(x,r) = \{x' \in \mathbb{R}^d : \|x-x'\|_2 \le \epsilon\}$ the ball centered in x of radius r. There exist methods that given first-order oracle access, only need $O(d\log 1/\epsilon)$ queries, and this query complexity is worst-case optimal (Nemirovsky et al. [28]) when $\epsilon \ll 1/\sqrt{d}$. Known methods achieving the optimal $O(d\log 1/\epsilon)$ query complexity fall in the broad class of cutting-plane methods that build upon the well-known ellipsoid method (Shor [35], Yudin and Nemirovskii [43]), which uses $O(d^2\log 1/\epsilon)$ queries. These include the inscribed ellipsoid (Nesterov [29], Tarasov [38]), volumetric center or Vaidya's method (Atkinson and Vaidya [2], Vaidya [39]), approximate center of mass via sampling techniques (Bertsimas and Vempala [5], Levin [19]), and recent improvements (Jiang et al. [16], Lee et al. [18]). Unfortunately, all of these methods suffer from at least $\Omega(d^3\log 1/\epsilon)$ time complexity, and they further require storing all subgradients or at least an ellipsoid in \mathbb{R}^d and therefore, at least $\Omega(d^2\log 1/\epsilon)$ bits of memory. These limitations are prohibitive for large-scale optimization; hence, cutting-plane methods are viewed as rather impractical and less frequently used for high-dimensional applications. On the other hand, the simplest, perhaps most commonly used and practical gradient descent requires $O(1/\epsilon^2)$ queries, which is not optimal for $\epsilon \ll 1/\sqrt{d}$, but only needs O(d) time per query and $O(d\log 1/\epsilon)$ memory.

A natural question is whether one can preserve the optimal query lower bounds from cutting-plane methods with simpler methods: for instance, inspired by gradient descent techniques. Such hope is largely motivated by the fact that in many different theoretical settings, cutting-plane methods have achieved state-of-the-art run times, including semidefinite programming (Anstreicher [1], Lee et al. [18]), submodular optimization (Grötschel

^{*}Corresponding author

et al. [13], Jiang [14], Lee et al. [18], McCormick [23]), or equilibrium computation (Jiang and Leyton-Brown [15], Papadimitriou and Roughgarden [32]). Toward this goal, Woodworth and Srebro [42] first posed this question in terms of query complexity/memory trade-off. Given a certain number of bits of memory, which query complexity is achievable? Although cutting-plane methods require $\Omega(d^2\log 1/\epsilon)$ memory, gradient descent only requires storing one vector, and as a result, it uses $O(d\log 1/\epsilon)$ memory, which is information-theoretically optimal (Woodworth and Srebro [42]); $\Omega(d\log 1/\epsilon)$ bits of memory are already required just to represent the answer to the optimization problem. Understanding this trade-off could pave the way for the design of more efficient methods in convex optimization.

The first result in this direction was provided in Marsden et al. [22], where they showed that it is impossible to be both optimal in query complexity and in memory. Specifically, they proved that any potentially randomized algorithm that uses at most $d^{1.25-\delta}$ memory must make at least $\tilde{\Omega}(d^{1+4/3\delta})$ queries for all $\delta \in [0,1/4]$. This implies that a superlinear amount of memory $d^{1.25}$ is required to achieve the optimal rate of convergence (that is achieved by algorithms using more than quadratic memory). However, this leaves open the fundamental question of whether one can improve over the memory of cutting-plane methods while keeping optimal query complexity.

Question (Conference on Learning Theory 2019 (Woodworth and Srebro [42])). Is it possible for a first-order algorithm that uses at most $O(d^{2-\delta})$ bits of memory to achieve query complexity $\tilde{O}(d \operatorname{polylog} 1/\epsilon)$ when $d = \Omega(\log^c 1/\epsilon)$ but $d = o(1/\epsilon^c)$ for all c > 0?

In this paper, building upon the techniques introduced in Marsden et al. [22], we provide a negative answer to this question; quadratic memory is necessary to achieve the optimal query complexity with deterministic algorithms. As a result, cutting-plane methods, including the standard center-of-mass algorithm, are Pareto optimal up to logarithmic factors within the query complexity/memory trade-off. Our main result for convex optimization is the following.

Theorem 1. For $\epsilon = 1/d^4$ and any $\delta \in [0,1]$, a deterministic first-order algorithm guaranteed to minimize 1-Lipschitz convex functions over the unit ball with ϵ accuracy uses at least $d^{2-\delta}$ bits or makes $\tilde{\Omega}(d^{1+\delta/3})$ queries.

A key component of cutting-plane methods is that they merely rely on the subgradient information at each query to restrict the search space. As a result, these can be used to solve the larger class of feasibility problems that are essential in mathematical programming and optimization. In a feasibility problem, one aims to find an ϵ -approximation of an unknown vector x^* and has access to a separation oracle. For any query x, the separation oracle either returns a separating hyperplane g from x to $B_d(x^*, \epsilon)$ —such that $\langle g, x-z \rangle > 0$ for any $z \in B_d(x^*, \epsilon)$ —or signals that $||x-x^*|| \le \epsilon$. This class of problems is broader than convex optimization because the negative subgradient always provides a separating hyperplane from a suboptimal query to the optimal set. Hence, feasibility and convex minimization problems are closely related, and it is often the case that obtaining query lower bounds for the feasibility problem simplifies the analysis while still providing key insights for the more restrictive convex optimization problem (Nemirovsky et al. [28], Nesterov [30]).

As a result, a similar fundamental question is to understand the query complexity/memory trade-off for the feasibility problem. As noted above, any lower bound for convex optimization yields the same lower bound for the feasibility problem. Here, we can significantly improve over the previous trade-off.

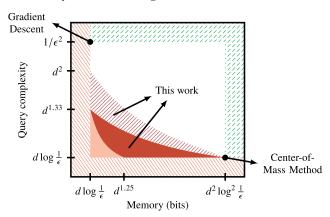
Theorem 2. For $\epsilon = 1/(48d^2\sqrt{d})$ and any $\delta \in [0,1]$, a deterministic algorithm guaranteed to solve the feasibility problem over the unit ball with ϵ accuracy uses at least $d^{2-\delta}$ bits of memory or makes at least $\tilde{\Omega}(d^{1+\delta})$ queries.

1.1. Literature Review

Recently, there has been a series of studies exploring the trade-offs between sample complexity and memory constraints for learning problems, such as linear regression (Sharan et al. [34], Steinhardt and Duchi [36]), principal component analysis (Mitliagkas et al. [24]), learning under the statistical query model (Steinhardt et al. [37]), and other general learning problems (Beame et al. [4], Brown et al. [7], Brown et al. [8], Garg et al. [12], Kol et al. [17], Moshkovitz and Moshkovitz [26]).

For parity problems that meet certain spectral (mixing) requirements, Raz [33] first proved by a computation tree argument that an exponential number of random samples is needed if the memory is subquadratic. Similar trade-offs have been obtained when the learning problem satisfies other types of properties (Beame et al. [4], Garg et al. [12], Kol et al. [17], Moshkovitz and Moshkovitz [25], Moshkovitz and Moshkovitz [26]). It should be noted that all of the above-mentioned results hold for learning problems over finite fields (i.e., the concept classes are finite). For continuous problems, Sharan et al. [34] was the first to apply the framework of Raz [33] and showed a sample complexity lower bound for memory-constrained linear regression.

Figure 1. (Color online) Trade-offs between available memory and first-order oracle complexity for minimizing 1-Lipschitz convex functions over the unit ball (adapted from Marsden et al. [22] and Woodworth and Srebro [40]). The bottom-left dashed "L"-shaped region (top-right dashed inverted "L"-shaped region, respectively) corresponds to historical information-theoretic lower bounds (upper bounds, respectively) on the memory and query complexity. The light bottom-left solid region corresponds to the recent lower-bound trade-off from Marsden et al. [22], which holds for randomized algorithms. In our work, we show that the dark solid region is not achievable for any deterministic algorithms. For the feasibility problem, we also show that the dark dashed region is not achievable either for any deterministic algorithms.



In contrast to learning with random samples, there is limited understanding of the memory-constrained optimization and feasibility problem. Nemirovsky et al. [28] demonstrated that in the absence of memory constraints, finding an ϵ -approximate solution for Lipschitz convex functions requires $\Omega(d\log 1/\epsilon)$ queries, which can be achieved by the center-of-mass method using $O(d^2\log^2 1/\epsilon)$ bits of memory. At the other extreme, gradient descent needs $\Omega(1/\epsilon^2)$ queries but only $O(d\log 1/\epsilon)$ bits of memory, the minimum memory needed to represent a solution. These two extreme cases are represented by the dashed pink "impossible region" and the dashed green "achievable region" in Figure 1. Since then, Marsden et al. [22] showed that there is a trade-off between memory and query for convex optimization; it is impossible to be both optimal in query complexity and memory. Their lower bound is represented by the solid pink "impossible region" in Figure 1. In this paper, we significantly improve these results to match the quadratic upper bound of cutting-plane methods. Additionally, there has been recent progress in the study of query complexity for randomized algorithms (Woodworth and Srebro [40], Woodworth and Srebro [41]).

On the algorithmic side, the aforementioned methods that achieve O(poly(d)) query complexity (Atkinson and Vaidy [2], Bertsimas and Vempala [5], Jiang et al. [16], Lee et al. [18], Levin [19], Nesterov [29], Shor [35], Tarasov [38], Vaidya [39], Yudin and Nemirovskii [43]) all require at least $\Omega(d^2\log 1/\epsilon)$ bits of memory. There is also significant literature on memory-efficient optimization algorithms, such as the limited-memory Broyden–Fletcher–Goldfarb–Shanno (BFGS) algorithm (Liu and Nocedal [21], Nocedal [31]). However, the convergence behavior for even the original BFGS on nonsmooth convex objectives is still a challenging, open question (Lewis and Overton [20]).

1.1.1. Comparison with Marsden et al. [22]. Our proof techniques build upon those introduced in Marsden et al. [22]. We follow the proof strategy that they introduced to derive lower bounds for the memory/query complexity. Below, we delineate which ideas and techniques are borrowed from Marsden et al. [22] and which are the novel elements that we introduce. Details on these proof elements are given in Section 2.4.

First, Marsden et al. [22] define a class of difficult functions for convex optimization of the following form:

$$F(\mathbf{x}) = \max \left\{ \|A\mathbf{x}\|_{\infty} - \eta_0, \eta_0 \left(\max_{i \le N} \mathbf{v}_i^{\mathsf{T}} \mathbf{x} - i \gamma \right) \right\},\tag{1}$$

where $A \sim \mathcal{U}(\{\pm 1\}^{d/2\times d})$ is a matrix with ± 1 entries sampled uniformly, and $v_i \sim \mathcal{U}(d^{-1/2}\{\pm 1\}^d)$ are sampled independently, uniformly within the rescaled hypercube. To give intuition on this class, the term $||Ax||_{\infty} - \eta_0$ acts as a barrier; to observe subgradients from the other term, one needs to use queries x that are approximately within the null space of A. The second term $\max_{i \le N} v_i^{\mathsf{T}} x - i \gamma$ is the "Nemirovski" function, which was used in previous works (Balkanski and Singer [3], Bubeck et al. [9], Nemirovski [27]) to obtain lower bounds in parallel convex

optimization. At a high level, the limitation in the lower bounds from Marsden et al. [22] comes from the fact that one is limited in the number N of vectors v_1, \ldots, v_N that can be used in the Nemirovski function. To resolve this issue, we introduce adaptivity within the choice of a modified Nemirovski function. At a high level, we choose the vectors v_1, \ldots, v_N depending on the queries of the algorithm, which allows us to fit in more terms. In turn, this allows us to improve the lower bounds.

As a second step, Marsden et al. [22] relate the optimization problem on the defined class of functions to an orthogonal vector game. In this game, the goal is to find vectors that are approximately orthogonal to a matrix A with access to row queries of A. The argument is as follows; because of the barrier term $||Ax||_{\infty} - \eta_0$, optimizing the Nemirovski function requires exploring independent directions of the null space of A, which is performed at informative queries. With our new class of functions, we can adapt this logic. However, the adaptivity in the vectors v_i provides information to the learner on A in addition to the queried rows of A. We, therefore, need to modify the game by introducing an orthogonal vector game with hints, where hints encapsulate this extra information.

For the last step, Marsden et al. [22] give an information-theoretic argument to provide a query complexity lower bound on the defined orthogonal vector game. Following the same structure, we show that a similar argument holds for our modified game. The main added difficulty resides in bounding the information leakage from the hints, and we show that these provide no more information than the memory itself.

As a last remark, the lower bounds provided in Marsden et al. [22] hold for randomized algorithms, whereas the adaptivity of our procedure only applies to deterministic algorithms.

1.2. Outline of the Paper

Our main results for the trade-off between memory and query complexity for optimization and the feasibility problem have been presented in Section 1 (Theorems 1 and 2). In Section 2, we formally define memory-constrained algorithms and provide a brief overview of our proof techniques and contributions. Our proofs for convex optimization are given in Section 3. We introduce the *optimization procedure*, which adaptively constructs a hard family of functions; additionally, we provide a reduction to this hard family from an *orthogonal vector game with hints*, and we show a memory-sample trade-off (Proposition 5) for the game, which completes the proof of Theorem 1. Last, in Section 4, we consider the feasibility problem and with a similar methodology, prove Theorem 2.

2. Formal Setup and Overview of Techniques

Standard results in oracle complexity give the minimal number of queries for algorithms to solve a given problem. However, this does not account for possible restrictions on the memory available to the algorithm. In this paper, we are interested in the trade-off between memory and query complexity for both convex optimization and the feasibility problem. Our results apply to a large class of *memory-constrained* algorithms. We give below a general definition of the memory constraint for algorithms with access to an oracle $\mathcal{O}: \mathcal{S} \to \mathcal{R}$ taking as input a query $q \in \mathcal{S}$ and returning as response $\mathcal{O}(q) \in \mathcal{R}$.

Definition 1 (*M*-Bit Memory-Constrained Deterministic Algorithm). Let $\mathcal{O}: \mathcal{S} \to \mathcal{R}$ be an oracle. An *M*-bit memory-constrained deterministic algorithm is specified by a query function $\psi_{query}: \{0,1\}^M \to \mathcal{S}$ and an update function $\psi_{update}: \{0,1\}^M \times \mathcal{S} \times \mathcal{R} \to \{0,1\}^M$. The algorithm starts with the memory state Memory₀ = 0^M and iteratively makes queries to the oracle. At iteration t, it makes the query $q_t = \psi_{query}(\mathsf{Memory}_{t-1})$ to the oracle, receives the response $r_t = \mathcal{O}(q_t)$, and then, updates its memory Memory_t = $\psi_{update}(\mathsf{Memory}_{t-1}, q_t, r_t)$.

The algorithm can stop making queries at any iteration, and the last query is its final output. Notice that the memory constraint applies only between each query but not for internal computations (i.e., the computation of the update ψ_{update} and the query ψ_{query} can potentially use unlimited memory). This is a rather weak memory constraint on the algorithm; a fortiori, our negative results also apply to stronger notions of memory-constrained algorithms. In Definition 1, we ask the query and update functions to be time invariant. In our context, this is without loss of generality; any M-bit algorithm using T queries with time-dependent query and update functions (Marsden et al. [22], Woodworth and Srebro [42]) can be turned into an $(M + \lceil \log T \rceil)$ -bit time-invariant algorithm by storing the iteration number t as part of the memory. The query lower bounds that we provide are at most $T \leq poly(d)$. Hence, the additional $\log T = O(\log d)$ bits to the memory size M do not affect our main results: Theorems 1 and 2.

In this paper, we use the above-described framework to study the interplay between query complexity and memory for two fundamental problems in optimization and machine learning.

2.1. Convex Optimization

We first consider convex optimization, in which one aims to minimize a 1-Lipschitz convex function $f: \mathbb{R}^d \to \mathbb{R}$ over the unit ball $B_d(0,1) \subset \mathbb{R}^d$. The goal is to output a point $\tilde{x} \in B_d(0,1)$ such that $f(\tilde{x}) \leq \min_{x \in B_d(0,1)} f(x) + \epsilon$, referred to as ϵ -approximate points. The optimization algorithm has access to a first-order oracle $\mathcal{O}_{CO}: \mathbb{R}^d \to \mathbb{R} \times \mathbb{R}^d$, which for any query x, returns the couple $(f(x), \partial f(x))$, where $\partial f(x)$ is a subgradient of f at the query point x.

Remark 1. The above requirement for ϵ -approximate optimality is weaker than asking to find a point that is at distance ϵ from arg $\min_{x \in B_d(0,1)} f(x)$ (for 1-Lipschitz convex functions). As a result, our lower bounds for ϵ -approximate optimality hold a fortiori for the problem where one aims to find a point at a distance at most ϵ from the solution set.

2.2. Feasibility Problem

Second, we consider the trade-off between memory and query complexity for the feasibility problem, where the goal is to find an element $\tilde{x} \in Q$ for a convex set $Q \subset B_d(0,1)$. Instead of a first-order oracle, the algorithm has access to a separation oracle $\mathcal{O}_F : \mathbb{R}^d \to \{\text{Success}\} \cup \mathbb{R}^d$. For any query $x \in \mathbb{R}^d$, the separation oracle either returns Success, reporting that $x \in Q$, or provides a separating vector $g \in \mathbb{R}^d$: that is, such that for all $x' \in Q$,

$$\langle g, x - x' \rangle > 0.$$

We say that an algorithm solves the feasibility problem with accuracy $\epsilon > 0$ if it can solve any feasibility problem for which the successful set contains a ball of radius ϵ (i.e., such that there exists $x^* \in B_d(0,1)$ satisfying $B_d(x^*,\epsilon) \subset Q$).

The feasibility problem is at least as hard as convex optimization in the following sense; an algorithm that solves the feasibility problem with accuracy ϵ/L can be used to solve L-Lipschitz convex optimization problems by feeding the subgradients from first-order queries to the algorithm as separating hyperplanes. Alternatively, from any 1-Lipschitz function f, one can derive a feasibility problem, where the feasibility set is $Q = \{x \in B_d(0,1), f(x) \le f^* + \epsilon\}$ and the separation oracle at $x \notin Q$ is a subgradient $\partial f(x)$ at x.

2.3. Overview of the Proof in Marsden et al. [22]

To ease the presentation, we first give an overview of the proof techniques from Marsden et al. [22], which we build upon. We recall that the family of functions that they use is given in Equation (1). The first term $||Ax||_{\infty} - \eta_0$ acts as a barrier term; to observe subgradients from the other term, one needs the query x to satisfy $||Ax||_{\infty} \le 2\eta_0$. These are called *informative queries*. They must lie approximately in the orthogonal space to the lines of A; that is, they approximately belong to the null space of A denoted Ker(A). Note that function F is designed so that intuitively, its minimum is given by the second term. Hence, an optimization algorithm needs to make informative queries in order to optimize F.

The second term $\max_{i \in [N]} v_i^\top x - i\gamma$ is referred to as a *Nemirovski* function. If γ is set appropriately to $\gamma = \Omega(\sqrt{\log d/d})$, an algorithm that optimizes this function must discover the subgradients v_1, \ldots, v_N in this exact order. In fact, for any $k \ge 1$, choosing $\gamma = \Omega(\sqrt{k \log d/d})$, they prove that (1) subgradients v_1, \ldots, v_N are discovered exactly in this order and that (2) any query that visits a new vector v_i does not lie close to the subspace formed by the last k informative vectors, a property known as *robust linear independence*. Indeed, for the last claim, from high-dimensional concentration, for a random unit vector v and a k-dimensional subspace E, $||P_E(v)|| = \Theta(\sqrt{k \log d/d})$.

As a result, at any point when optimizing F, in order to observe the next k subgradients from the Nemirovski function, one needs to make k informative queries that are robustly independent and close to Ker(A). The crux of the optimization difficulty is that the algorithm receives information about A only through the subgradients of $||Ax||_{\infty}$: that is, one row at a time. This motivates the definition of the following (simplified) game.

- 1. *Oracle*. Sample $A \sim \mathcal{U}(\{\pm 1\}^{d/2 \times d})$.
- 2. Player. Based on A, store an M-bit message Message.
- 3. *Player*. Using only Message (but not *A*), query some rows of *A*, and output vectors y_1, \ldots, y_k .
- 4. The player wins if the returned vectors are all approximately in Ker(A) and are robustly independent; that is, $||P_{Span(y_i,j< i)^{\perp}}(y_i)|| \ge \beta$ for all $i \in [k]$ for some fixed parameter β .

They show that to win this game, the player should either (1) make $\Omega(d)$ row queries or (2) use memory $M = \Omega(kd)$. Roughly speaking, their result shows that to find k robustly independent vectors roughly in Ker(A), either (1) we need to query all rows of A (once we know A, finding vectors in its null space is easy), or (2) we need to store these vectors directly in memory, which requires $\tilde{O}(kd)$ bits of memory. Setting $k \approx CM/d$ for some large constant C, where M is the bit memory of the algorithm, ensures that only the first scenario happens.

With these ingredients at hand, assuming that the algorithm needs to discover all N subgradients v_1, \ldots, v_N of the Nemirovski term, this gives a query lower bound of $\Omega(d) \times (N/k)$. Unfortunately, the maximum number of useful Nemirovski vectors N is bounded by the value of γ ; one needs $N \leq N_0 = \tilde{O}(\gamma^{-2/3})$. Beyond this value, for any $j > N_0$, we would have $v_j^\top x - j\gamma \leq \max_{i \in [N_0]} v_i^\top x - i\gamma$ for all $x \in B_d(0,1)$; hence, further terms are irrelevant to optimize F. This gives a final query lower bound of $\Omega(Nd/k) = \Omega(Nd^2/M) = \tilde{\Omega}((d^2/M)^{4/3})$ for M-bit memory algorithms.

2.4. Overview of the Proof Strategy and Innovations

Because the techniques for Theorems 1 and 2 are similar, we mostly focus on main ideas used to derive lower bounds for convex optimization. Although our proof borrows techniques from Marsden et al. [22], we introduce key innovations involving adaptivity to improve the lower bounds up to the maximum quadratic memory for deterministic algorithms—up to logarithmic factors. We recall, however, that the bounds in Marsden et al. [22] hold for randomized algorithms as well. In the proofs, we aim to optimize the dependence of the parameters in d. Constants, however, are not necessarily optimized.

As a road map, our proof has three main components (see Figure 2). We first show that solving the general memory-constrained *convex optimization* problem implies solving an *optimization procedure* (Proposition 1). Necessary properties on the optimization procedure are proven in Propositions 2 and 3. We then further relate the optimization procedure to an *orthogonal vector game with hints* (Proposition 4), on which we prove memory/query trade-offs in Proposition 5.

2.4.1. Main Limitations for Improving the Lower Bounds. As per the computations in Section 2.3, one of the main barriers to improving the lower bounds is the limit on the number N of Nemirovski vectors that can be used. Ideally, if one could ensure $N = \Omega(d)$, which is the maximum possible value, then this would directly give a lower-bound trade-off up to quadratic memory $O(d^2)$. Our adaptive construction uses a different form of functions, but roughly speaking, we will be able to ensure precisely $N = \tilde{\Omega}(d)$ for the feasibility problem. However, for the optimization problem, we will only be able to reach the value $N = k(d/k)^{1/3}$, which still provides a query lower bound of $\Omega(Nd/k) = \Omega(d(d/k)^{1/3}) = \tilde{\Omega}(d(d^2/M)^{1/3})$.

As a preview, given the bound $N_0 = O(\gamma^{-2/3})$, one of the goals of the adaptive construction is to decrease the value of γ from $\Omega(\sqrt{k \log d/d})$ to $O(\sqrt{\log d/d})$, which is the minimum value that still ensures the subgradients v_1, \ldots, v_N to be observed in this exact order. We also use a two-layer construction to further reduce the value of γ for the last layer, which we discuss below. We recall that in the construction of Marsden et al. [22], having $\gamma = \Omega(\sqrt{k \log d/d})$ was necessary for k informative queries to be robustly independent.

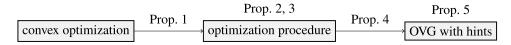
2.4.2. An Adaptive Optimization Procedure. Instead of using a fixed distribution of convex functions as a hard instance as in Equation (1), we construct the hard functions adaptively. To do so, we design an *optimization procedure*, which for any algorithm, constructs a hard family of convex functions adaptively on its queries from the following family of convex functions with appropriately chosen parameters η , γ_1 , γ_2 , p_{max} , l_p , $\delta > 0$:

$$F_{A,v}(x) = \max\left\{ ||Ax||_{\infty} - \eta, \eta v_0^{\top} x, \eta \left(\max_{p \le p_{max}, l \le l_p} v_{p,l}^{\top} x - p \gamma_1 - l \gamma_2 \right) \right\}.$$
 (2)

We take $A \sim \mathcal{U}(\{\pm 1\}^{n\times d})$ and $v_0 \sim \mathcal{U}(\mathcal{D}_\delta)$ uniformly sampled in the beginning, where $\mathcal{D}_\delta \subset \mathcal{S}^{d-1}$ is a (finite) discretization of the sphere. As in Equation (1), these functions include the barrier term $||Ax||_{\infty} - \eta$, and queries x that satisfy $||Ax||_{\infty} \leq 2\eta$ are called informative; these lie approximately in Ker(A). The second term $\eta v_0^{\mathsf{T}} x$ is used to ensure that solutions with low objective (in particular, with the objective at most $\eta \gamma_1/2$) have norm bounded away from zero. As a result, these informative queries, once renormalized, will still belong approximately to the null space of A denoted Ker(A).

The main novelty in the construction is captured by the third term, which is constructed adaptively along the optimization process. This construction proceeds by periods $p = 1, 2, ..., p_{max}$ designed so that during each period $p \in [p_{max}]$, the algorithm is forced to visit a subspace of Ker(A) of fixed dimension $k = \tilde{\Omega}(M/d)$. Here, k is a parameter that plays the same role as in Marsden et al. [22]; assuming that the algorithm needs to visit a subspace of dimension k, then it should make at least k queries approximately in Ker(A) that are robustly linearly

Figure 2. General proof structure. OVG, orthogonal vector game; Prop., proposition.



independent. The hope is that because $k = \tilde{\Omega}(M/d)$ (the algorithm cannot store these queries directly in memory), this requires making $\tilde{\Omega}(d)$ queries to the gradient oracle, yielding a final query lower bound of $\tilde{\Omega}(p_{max}d)$.

For each period p, to ensure that the algorithm visits a subspace of Ker(A) of dimension k, we iteratively construct vectors $v_{p,1}, \ldots v_{p,l_v}$ as follows. Suppose that at the beginning of a step of period p, one has defined vectors $v_{p,1}, \ldots, v_{p,l_v}$.

- The procedure first evaluates the explored subspace of the algorithm during this period. More precisely, the procedure keeps track of *exploratory* queries $x_{i_{p,1}}, \ldots, x_{i_{p,r}}$ during period p up to the current step. The exploratory subspace is then $Span(x_{i_{p,1}}, \ldots, x_{i_{p,r}})$.
- If a query with a sufficiently low objective is queried, we sample a new vector $v_{p,l+1}$, which is approximately orthogonal to the exploratory subspace. The corresponding new term in the objective is $v_{p,l+1}^{\top}x p\gamma_1 (l+1)\gamma_2$. Once this new term is added to the objective, the algorithm is constrained to make queries with an additional component along the direction $-v_{p,l+1}$. Because this vector is approximately orthogonal to all previous queries, this forces the algorithm to query vectors linearly independent from all previous queries in period p. The period then ends once the dimension of the exploratory subspace reaches k, having defined l_p vectors $v_{p,1}, \ldots, v_{p,l_p}$. As discussed above, the exploratory subspace must increase dimension for any additional such vector. Thus, after $l_p \leq k$ vectors, period p ends (Lemma 2).

As a comparison with the family of functions from Equation (1), the third term of Equation (2) now plays the role of the Nemirovski function, and the total number of Nemirovski terms is intuitively $N \approx p_{max}k$ because each layer $p \in [p_{max}]$ constructs $l_p \leq k$ vectors $v_{p,1}, \ldots, v_{p,l_n}$.

2.4.3. Benefits of Adaptivity. We now expand on how the adaptive terms allow for improving the lower bound of Marsden et al. [22] to match the quadratic upper bound of cutting-plane methods. As we mentioned above, one of the limitations in the functions of the form Equation (1) comes from the fact that the offset in the Nemirovski function is $\gamma = \Omega(\sqrt{k \log d/d})$. This offset was necessary to ensure that with high probability, (1) subgradients v_1, \ldots, v_N are discovered exactly in this order—in fact, this is also ensured whenever $\gamma = \Omega(\sqrt{\log d/d})$ —and (2) any query that visits a new vector v_i must not lie in the subspace formed by the last k last informative vectors. In our procedure, however, this is not necessary anymore because during each period $p \in [p_{max}]$, a k-dimensional subspace of Ker(A) is forced to be explored; that is, property (2) is already satisfied. Hence, we only need to ensure property (1). More precisely, we check that the optimization procedure is equivalent to running the optimization algorithm with the final constructed function. This is the goal of Proposition 1, and we provide the main ideas below.

We need to ensure that the algorithm first observes the subgradients of period 1 (that is, $v_{1,1}, \ldots, v_{1,l_1}$) and then those of period 2 until those of period p_{max} . This is the purpose of the offset γ_1 , which can, therefore, be taken as $\gamma_1 = O(\sqrt{\log d/d})$.

Within each period p, we also need to ensure that the subgradients $v_{p,1}, \ldots, v_{p,l_p}$ would be discovered in that order when running the optimization algorithm with the final constructed function. For this second layer, we will be able to further reduce the value of the offset γ_2 . The vectors $v_{p,l}$ for $l \in [l_p]$ are constructed so that they are approximately orthogonal to any query x that was previously made during period p. Hence, we will be able to show that

$$\max_{l \le l' \le l_p} \mathbf{v}_{p,l'}^{\top} \mathbf{x} - p \gamma_1 - l' \gamma_2 \le -p \gamma_1 - (l-1) \gamma_2 - \frac{\gamma_2}{2}. \tag{3}$$

This gives an offset $-\gamma_2/2$ compared with the previous terms $\max_{l' \leq l-1} v_{p,l'}^\top x - p \gamma_1 - l' \gamma_2$, which in turn, ensures that such a query x could not have observed the vectors $v_{p,l'}$ for $l' \geq l$. In fact, we can take γ_2 as small as desired; taking $\gamma_2 = \tilde{O}(\gamma_1/d)$ is sufficient. Because there are at most $l_p \leq k$ terms per period, the total offset per period still satisfies $l_p \gamma_2 \leq k \gamma_2 \ll \gamma_1$. In words, because the vectors $v_{p,l}$ are constructed perpendicular to exploration spaces within each period, the offset needed within each period is negligible.

Now that the offset parameters γ_1 and γ_2 have been reduced, we can increase the number of useful Nemirovski terms. Formally, it remains to estimate the maximum number of periods p_{max} that we can fit in the construction. First, using standard arguments, we show (Proposition 2) that

$$\min_{x \in B_d(0,1)} \frac{F_{A,v}(x)}{\eta} \leq -\tilde{\Omega}\left(\frac{1}{\sqrt{N}}\right),\,$$

where $N = p_{max}k$ is roughly the number of Nemirovski terms. On the other hand, suppose that an algorithm does not complete all p_{max} periods; say it does not observe $v_{p,l}$. Then, if x_T is the output of the algorithm, by concentration bounds, we have

$$\boldsymbol{v}_{p,l}^{\top}\boldsymbol{x}_T - p\boldsymbol{\gamma}_1 - l\boldsymbol{\gamma}_2 \geq -(p+1)\boldsymbol{\gamma}_1 - l\boldsymbol{\gamma}_2 \geq -O(p_{max}\boldsymbol{\gamma}_1).$$

Note that if we choose $p_{max} = \tilde{\Theta}((d/k)^{1/3})$, we have $1/\sqrt{N} = 1/\sqrt{p_{max}k} \gg p_{max}\gamma_1$. Then, combining the three previous inequalities precisely shows that for $p_{max} = \tilde{\Theta}((d/k)^{1/3})$, an algorithm needs to complete p_{max} periods in order to find an approximate minimizer of $F_{A,v}$. Hence, the total number of Nemirovski terms is $\tilde{\Omega}(k(d/k)^{1/3})$ as desired. Full details are given in Proposition 3, which completes the reduction from the optimization procedure to convex optimization.

2.4.4. An Orthogonal Vector Game with Hints. A crucial part of the proof is to prove that with the constructed optimization procedure, at each period p, to find k exploratory queries approximately in Ker(A) and robustly independent, the algorithm needs to perform $\tilde{\Omega}(d)$ queries to the gradient oracle.

We link the optimization of the above-mentioned constructed functions with an orthogonal vector game with hints (Proposition 4). As in the game introduced by Marsden et al. [22] (see Section 2.3), the goal for the player is to find k linearly independent vectors approximatively in Ker(A). To do so, the player can access an M-bit message Message and make m queries to rows of A. We now give some brief intuition about their information-theoretic query lower bound. Suppose that $M \le ckd$ for a small constant c > 0. To win, the output vectors y_1, \ldots, y_k should be robustly independent, which intuitively implies that the algorithm needs to visit roughly k distinct dimensions of Ker(A). Each new dimension of Ker(A) must be (approximately) orthogonal to all lines of A. Hence, this provides additional mutual information O(k) for every line of A, including the d/2 - m lines that were not observed through queries. This extra information on A can only be explained by the message, which has M bits. Hence, $M \ge O(k)(d/2 - m)$. Setting the constant c > 0 appropriately, this shows that $m = \Omega(d)$.

In our case, the optimization procedure ensures that the algorithm needs to explore k dimensions of Ker(A) in each period. However, each query yields a response from the optimization oracle that can either be a line of A (corresponding to the term $||Ax||_{\infty} - \eta$ of Equation (2)) or v_0 (term $\eta v_0^{\top} x$ of Equation (2)) or be previously defined vectors $v_{p',l'}$. Now, because the vectors $v_{p',l'}$ have been constructed adaptively on the queries of the algorithm, which themselves may depend on lines of A, during a period p, responses $v_{p',l'}$ for p' < p are a source of information leakage for A from previous periods. As a result, the query lower bound on the game introduced by Marsden et al. [22] is not sufficient for our purposes. Instead, we introduce an orthogonal vector game with hints, where hints correspond exactly to these vectors $v_{p',l'}$ from previous periods. Informally, the game corresponds to a simulation of one of the periods of the optimization procedure; for each query x, the oracle returns the subgradient that would have been returned in the optimization procedure, up to minor details. The formal definition of the orthogonal vector game with hints is given in Game 1; we give here its general structure for intuition.

- 1. *Oracle*. Sample $A \sim \mathcal{U}(\{\pm 1\}^{d/4 \times d})$.
- 2. *Player*. Based on A, construct vectors v_1, \ldots, v_d according to a specific procedure that mimics the construction of vectors $v_{v,l}$ in the optimization procedure.
- 3. *Player*. Based on A, store an M-bit message Message, and submit a response function g, which takes values in $B_d(0,1)$ and outputs wither a row of A or a vector from v_1, \ldots, v_d .
- 4. *Player*. Using only Message (but not A), make some queries to the response function g, and output vectors y_1, \dots, y_k .
 - 5. The player wins if the returned vectors are approximately in Ker(A) and robustly independent.

2.4.5. Bounding the Information Leakage. The next step of the proof is to give lower bounds on the number of queries needed to solve the orthogonal vector game with hints (Proposition 5). The main difficulty is to bound the information leakage from these hints. We recall that hints are of the form $v_{p',l'}$, which have been constructed adaptively on the queries of the algorithm during period p'. In particular, these contain information on the lines of A queried during period p' < p, which may be complementary with those queried during period p. If this total information leakage through the hints yields a mutual information with Ker(A) significantly higher than that of the M bits of Message, obtained lower bounds cannot possibly reflect any trade-off with memory constraints. It is, therefore, essential to obtain information leakage at most $O(M) = \tilde{O}(dk)$.

To solve this issue, we introduce a discretization \mathcal{D}_{δ} of the unit sphere where the vectors $v_{p,l}$ take value. Next, we show that each individual vector $v_{p',l'}$ from previous periods can only provide information $\tilde{O}(k)$ on the matrix A. To have an intuition on this, note that for any (at most) k vectors x_1, \ldots, x_k , the volume of the subset of the unit sphere S^{d-1} of vectors approximately orthogonal to x_1, \ldots, x_k , say $S(x_1, \ldots, x_k) = \{y \in S^{d-1} : |y^{\top}x_i| \leq d^{-3}, i \leq k\}$, is $q_k = O(1/d^{3k})$. Hence, because the vector v is roughly taken uniformly at random within $\mathcal{D}_{\delta} \cap S(x_1, \ldots, x_k)$, we can

show that the mutual information of v with the initial vectors x_1, \ldots, x_k is at most $O(-\log q_k) = O(k\log d)$. As a result, even if m = d, the total information leakage through the vectors $v_{p',l'}$ from previous periods is at most $O(kd\log d)$. This is comparable with the information of Message; hence, the main information-theoretic argument can be conserved. The formal proof involves anticoncentration bounds on the distance of a random unit vector to a linear subspace of dimension k (Lemma 4) as well as a more involved discretization procedure than the one presented above. In summary, by introducing adaptive functions through the optimization procedure, we show that the same memory-sample trade-off holds for the orthogonal vector game with hints and the game without hints introduced in Marsden et al. [22], up to logarithmic factors.

3. Memory-Constrained Convex Optimization

To prove our results, we need to use discretizations of the unit sphere S^{d-1} . It will be convenient to ensure that the partitions induced by these discretizations have equal area, which can be done with the following lemma.

Lemma 1 (Feige and Schechtman [11, Lemma 21]). For any $0 < \delta < \pi/2$, the sphere S^{d-1} can be partitioned into $N(\delta) = (O(1)/\delta)^d$ equal volume cells, each of diameter at most δ .

We denote by $V_{\delta} = \{V_i(\delta), i \in [N(\delta)]\}$ the corresponding partition and consider a set of representatives $\mathcal{D}_{\delta} = \{b_i(\delta), i \in [N(\delta)]\} \subset S^{d-1}$ such that for all $i \in [N(\delta)]$, $b_i(\delta) \in V_i(\delta)$. With these notations, we can define the discretization function ϕ_{δ} as follows:

$$\phi_{\delta}(\mathbf{x}) = \mathbf{b}_i(\delta), \quad \mathbf{x} \in V_i(\delta).$$

3.1. Definition of the Difficult Class of Optimization Problems

In this section, we present the class of functions that we use to prove our lower bounds. Throughout the paper, we pose $n = \lceil d/4 \rceil$. We first define some useful functions. For any $A \in \mathbb{R}^{n \times d}$, we define g_A as follows:

$$g_A(x) = a_{i_{\min}}, \quad i_{\min} = \min\{i \in [n], |a_i^\top x| = ||Ax||_{\infty}\}.$$

With this function, we can define a subgradient function for $x \mapsto ||Ax||_{\infty}$:

$$\tilde{\mathbf{g}}_A(\mathbf{x}) = \epsilon \mathbf{g}_A(\mathbf{x}), \qquad \epsilon = sign(\mathbf{g}_A(\mathbf{x})^{\mathsf{T}}\mathbf{x}).$$

We are now ready to introduce the class of functions, which we use for our lower bounds. These are of the following form:

$$F_{A,v}(x) = \max \left\{ \|Ax\|_{\infty} - \eta, \eta v_0^{\top} x, \eta \left(\max_{p \leq p_{max}} \max_{l \leq l_p} v_{p,l}^{\top} x - p \gamma_1 - l \gamma_2 \right) \right\}.$$

Here, $A \in \{\pm 1\}^{n \times d}$ is a matrix. Also, v_0 and the terms $v_{p,l}$ are vectors in \mathbb{R}^d . More precisely, these vectors will lie in the discretization \mathcal{D}_δ for $\delta = 1/d^3$. We postpone the definition of p_{max} and l_p for $p \leq p_{max}$. Last, we use the following choice for the remaining parameters: $\eta = 2/d^3$, $\gamma_1 = 12\sqrt{\log d/d}$, and $\gamma_2 = \gamma_1/4d$. For convenience, we also define the functions

$$\begin{split} F_{A}(x) &= \max\{\|Ax\|_{\infty} - \eta, \eta v_0^{\top} x\} \\ F_{A,v,p,l}(x) &= \max\left\{\|Ax\|_{\infty} - \eta, \eta v_0^{\top} x, \eta \left(\max_{(p',l') \leq l_{ex}(p,l), l' \leq l_{p'}} v_{p',l'}^{\top} x - p' \gamma_1 - l' \gamma_2\right)\right\}, \end{split}$$

with the convention $F_{A,v,1,0} = F_A$. The functions $F_{A,v,p,l}$ will encapsulate the current state of the function to be minimized; it will be updated adaptively on the queries of the algorithm. We also define a subgradient function for $F_{A,v,p,l}$ by first favoring lines of A and then vectors from v in case of ties as follows:

$$\partial F_{A,v,p,l}(x) = \begin{cases} \tilde{g}_A(x_t) & \text{if } F_{A,v,l,p}(x) = ||Ax||_{\infty} - \eta, \\ \eta v_0 & \text{otherwise and if } F_{A,v,l,p}(x) = \eta v_0^\top x, \\ \eta v_{p,l} & \text{otherwise and if } (p,l) = \underset{(p',l') \leq_{lex}(p,l)}{\arg\max} v_{p',l'}^\top x - p' \gamma_1 - l' \gamma_2. \end{cases}$$

In the last case, ties are broken by lexicographic order. We define $\partial F_{A,v} = \partial F_{A,v,p_{max},l_{vmax}}$ similarly.

We consider a so-called *optimization procedure*, which will construct the sequence of vectors $v = (v_{p,l})$ adaptively on the responses of the considered algorithm. Throughout this section, we use a parameter $1 \le k \le d/3 - 1$ —which will be taken as $k = \Theta(M/d)$, where M is the memory of the algorithm—and let p_{max} be the largest number that satisfies the following constraint:

$$p_{max} \le \min\{(c_{d,1}d - 1)/k, c_{d,2}(d/k)^{1/3} - 1\},\tag{4}$$

where $c_{d,1} = 1/(90^2 \log^2 d)$ and $c_{d,2} = 1/(81 \log^{2/3} d)$.

Procedure 1 (The Optimization Procedure for Algorithm alg)

```
Input: d, k, p_{max}, algorithm alg
```

Part 1: Procedure to adaptively construct *v*

Sample $A \sim \mathcal{U}(\{\pm 1\}^{n \times d})$ and $v_0 \sim \mathcal{U}(\mathcal{D}_{\delta})$

2 Initialize the memory of alg to **0**, and let p = 1, r = l = 0

3 for $t \ge 1$, do

if $t > d^2$, **then** set (P, L) = (p, l), and break the **for** loop;

5 Run *alg* with current memory to obtain a query x_t

if $F_A(x) > \eta$, then //Noninformative query

return ($||Ax_t||_{\infty} - \eta, \tilde{g}_A(x_t)$) as response to *alg* 7

8 else //Informative query

if $r \le k-1$ and $F_{A,v,p,l}(x_t) \le -\eta \gamma_1/2$ and $\|P_{Span(x_{i_n,r},r'\le r)^{\perp}}(x_t)\|/\|x_t\| \ge \frac{\gamma_2}{4}$, then 9 10

Set $i_{v,r+1} = t$ and increment $r \leftarrow r + 1$

if $F_{A,v,p,l}(x_t) < -\eta(p\gamma_1 + l\gamma_2 + \gamma_2/2)$ and r < k, then 11

Compute Gram–Schmidt decomposition $b_{p,1},\ldots,b_{p,r}$ of $x_{i_{p,1}},\ldots,x_{i_{p,r}}$ Sample $y_{p,l+1}$ uniformly on $S^{d-1}\cap\{z\in\mathbb{R}^d:|b_{p,r'}^\top z|\leq d^{-3},\ \forall r'\leq r\}$ 13

Define $v_{p,l+1} = \phi_{\delta}(y_{p,l+1})$ and increment $l \leftarrow l+1$

15 else if $F_{A,v,p,l}(x_t) < -\eta(p\gamma_1 + l\gamma_2 + \gamma_2/2)$ and $p+1 \le p_{max}$, then

Set $l_p = l$ and $i_{p+1,1} = t$

17 Compute the Gram–Schmidt decomposition $b_{p+1,1}$ of $x_{i_{p+1,1}}$

Sample $y_{p+1,1}$ uniformly on $\mathcal{S}^{d-1} \cap \{z \in \mathbb{R}^d : |b_{p+1,1}^\top z| \le d^{-3}\}$ 18

Define $v_{p+1,1} = \phi_{\delta}(y_{p+1,1})$, increment $p \leftarrow p+1$, and reset l=r=119

else if $F_{A,v,p,l}(x_t)<-\eta(p\gamma_1+l\gamma_2+\gamma_2/2)$, then //End of the construction 20

21 Set $l_{p_{max}} = l$, $i_{p_{max}+1,1} = t$

Set $(P, L) = (p_{max}, l)$, and break the **for** loop 22

23 **return** $(F_{A,v,p,l}(x_t), \partial F_{A,v,p,l}(x_t))$ as response to *alg*

24 end

12

14

16

Part 2: Procedure once v, P, L are constructed

25 **for** $t' \ge t$, **do return** $(F_{A,v,P,L}(x_{t'}), \partial F_{A,v,P,L}(x_{t'}))$ as response to the query $x_{t'}$

The optimization procedure is described in Procedure 1. First, we sample independently $A \sim \mathcal{U}(\{\pm 1\}^{n \times d})$ and $v_0 \sim \mathcal{U}(\mathcal{D}_\delta)$. The matrix A and vector v_0 are then fixed for the rest of the learning procedure. Next, we describe the adaptive procedure to return subgradients. It proceeds by periods until p_{max} periods are completed unless the total number of iterations reaches d^2 , in which case the construction procedure ends as well. First, we say that a query is informative if $F_A(x) \le \eta$. The procedure proceeds by periods $p \in [p_{max}]$ and in each period, constructs the vectors $v_{p,1}, \ldots, v_{p,k}$ iteratively. We are now ready to describe the procedure at time t when the new query x_t is queried. Let $p \ge 1$ be the index of the current period and $v_{p,1}, \dots, v_{p,l}$ be the vectors of this period constructed so far; the first period is p = 1, and we allow l = 0 here. As will be seen in the construction, we always have $l \ge 1$ except at the very beginning, for which we use the notation $F_{A,v,1,0} = F_A$. Together with these vectors, the oracle keeps in memory indices $i_{p,1}, \dots, i_{p,r}$ with $r \leq k$ of *exploratory* queries. The constructed vectors from previous periods are $v_{p',l'}$ for p' < p and $l' \le l_{p'}$.

(Case 1) If x_t is not informative (i.e., $F_A(x) > \eta$), then the procedure returns ($||Ax_t||_{\infty} - \eta, \tilde{g}_A(x_t)$).

(Case 2) Otherwise, we follow the next steps. If $r \le k - 1$,

$$F_{A,v,p,l}(x_t) \le -\frac{\eta \gamma_1}{2}$$
, and $\frac{\|P_{Span(x_{i_{p,r'}},r' \le r)^{\perp}}(x_t)\|}{\|x_t\|} \ge \frac{\gamma_2}{4}$,

we set $i_{p,r+1} = t$ and increment r. In this case, we say that x_t is exploratory.

(Case 2a) Recalling that $F_{A,v,p,l}$ is constructed so far, if $F_{A,v,p,l}(x_t) \ge \eta(-p\gamma_1 - l\gamma_2 - \gamma_2/2)$, we do not do anything.

(Case 2b) Otherwise and if r < k, let $b_{p,1}, \ldots, b_{p,r}$ be the result from the Gram–Schmidt decomposition of $x_{i_{p,1}}, \ldots, x_{i_{p,r}}$. Then, let $y_{p,l+1}$ be a sample of the distribution obtained by the uniform distribution $y_{p,l+1} \sim \mathcal{U}\left(S^{d-1} \cap \{z \in \mathbb{R}^d : |b_{p,r'}^{\top}z| \leq d^{-3}, \ \forall r' \leq r\}\right)$. We then pose $v_{p,l+1} = \phi_{\delta}(y_{p,l+1})$. Having defined this new vector, we increment l.

(Case 2c) Otherwise, if r = k, this ends period p. We write the total number of vectors defined during period p as $l_p := l$. If $p+1 \le p_{max}$, period p+1 starts from $t = i_{p+1,1}$. Similarly to above, let $b_{p+1,1}$ be the result of the Gram–Schmidt procedure on $x_{p+1,1}$, and we sample $y_{p+1,1}$ according to a uniform distribution $y_{p+1,1} \sim \mathcal{U}\left(S^{d-1} \cap \{z \in \mathbb{R}^d : |b_{p+1,1}^\top z| \le d^{-3}\}\right)$. Then, we pose $v_{p+1,1} = \phi_\delta(y_{p+1,1})$. We can then increment p and reset l = r = 1.

After these steps, with the current values of p and l, we return $(F_{A,v,p,l}(x_t), \partial F_{A,v,l,p}(x_t))$.

If we finished the last period $p = p_{max}$ or if we reached a total number of iterations d^2 , the construction phase of the function ends. In both cases, let us denote by P, L the last defined period and vector $v_{P,L}$. In particular, we have $p \le p_{max}$. From now on, the final function to optimize is $F_{A,v,P,L}$, and the oracle is a standard first-order oracle for this function using the subgradient function $\partial F_{A,v,P,L}$.

We will relate this procedure to the standard convex optimization problem and prove query lower bounds under memory constraints for this procedure. Before doing so, we formally define what we mean by solving this optimization procedure.

Definition 2. Let alg be an algorithm for convex optimization. We say that an algorithm alg is successful for the optimization procedure with probability $q \in [0,1]$ and accuracy $\epsilon > 0$ if taking $A \sim \mathcal{U}(\{\pm 1\}^{n \times d})$, running alg with the responses given by the procedure, and denoting by $x^*(alg)$ the final answer returned by alg, with probability at least q over the randomness of A, v_0 and of the procedure, one has

$$F_{A,v,P,L}(x^*(alg)) \le \min_{x \in B_d(0,1)} F_{A,v,P,L}(x) + \epsilon.$$

3.2. Properties and Validity of the Optimization Procedure

We begin this section with a simple lemma showing that during each period p, at most $l_p \le k$ vectors $v_{p,1}, \ldots, v_{p,l_p}$ are constructed.

Lemma 2. At any time of the construction procedure, $l \le r$. In particular, because $r \le k$, we have $l_p \le k$ for all periods $p \le p_{max}$.

Proof. Fix a period p. We prove this by induction. The claim is satisfied for any l=1 when $p\geq 2$ because in this case, at the first time $t=i_{p,1}$ of the period p, we also construct the first vector $v_{p,1}$. For p=1, note that the first informative query t that falls in case (2b) or case (2c) is exploratory. Indeed, in these cases, we have $F_{A,v,1,0}(x_t) < \eta(-\gamma_1-\gamma_2/2) \leq -\eta\gamma_1/2$, and the second criterion for an exploratory query is immediate $\|P_{Span(x_{i_1,r},r',\leq 0)}(x_t)\| = 0$ because no indices $i_{1,r'}$ have been defined yet.

We now suppose that the claim holds for $l-1 \ge 1$. Let $t_{p,l}$ be the time when $v_{p,l}$ is constructed and $i_{p,1}, \ldots, i_{p,r}$ be the indices constructed until the beginning of iteration $t_{p,l}$. If a new index $i_{p,r}$ was constructed in times $(t_{p,l-1},t_{p,l})$, then the claim holds immediately. Suppose that this is not the case. Note that $t_{p,l}$ falls in case (2b), which means in particular that

$$\eta(v_{p,l-1}^{\top}x_{t_{p,l}} - p\gamma_1 - (l-1)\gamma_2) \leq F_{A,v,p,l-1}(x_{t_{p,l}}) < \eta(-p\gamma_1 - (l-1)\gamma_2 - \gamma_2/2).$$

As a result,

$$|y_{p,l-1}^{\top} x_{t_{p,l}}| \ge |v_{p,l-1}^{\top} x_{t_{p,l}}| - \delta > \frac{\gamma_2}{2} - \delta.$$

Next, when $r \ge l-1$ is the number of indices constructed so far, we decompose $y_{p,l-1} = \alpha_1 b_{p,1} + \cdots + \alpha_r b_{p,r} + \tilde{y}_{p,l-1}$, where $\tilde{y}_{p,l-1} \in Span(x_{i_p,r},r' \le r)^{\perp}$. Now, by construction of $y_{p,l-1}$, one has $|\alpha_{r'}| \le d^{-3}$ for all $r' \le r$. Thus,

$$\|\tilde{y}_{p,l-1} - y_{p,l-1}\| \leq \frac{\sqrt{r}}{d^3} \leq \frac{1}{d^2 \sqrt{d}}.$$

Therefore,

$$||P_{Span(x_{i_{p,r'}}, r' \leq r)^{\perp}}(x_{t_{p,l}})|| \geq ||\tilde{y}_{p,l-1}^{\top} x_{t_{p,l}}|| \geq ||y_{p,l-1}^{\top} x_{t_{p,l}}|| - \frac{1}{d^2 \sqrt{d}} > \frac{\gamma_2}{2} - \frac{1}{d^2 \sqrt{d}} - \delta \geq \frac{\gamma_2}{4}.$$

As a result, $t_{v,l}$ is exploratory; hence, $i_{v,r+1} = t_{v,l}$. This ends the proof of the recursion and the lemma. \Box

We recall that P and L denote the last defined period and vector $v_{P,L}$. From Lemma 2, we have in particular $P \le p_{max}$ and $L \le k$. In the next result, we show that with high probability, the returned values and vectors returned by the above procedure are consistent with a first-order oracle for minimizing the function $F_{A,v,P,L}$.

Proposition 1. Let $A \in \{\pm 1\}^{n \times d}$ and $v_0 \in \mathcal{D}_{\delta}$. On an event \mathcal{E} of probability at least $1 - C\sqrt{\log d}/d^2$ on the randomness of the procedure for some universal constant C > 0, all responses of the optimization procedure are consistent with a first-order oracle for the function $F_{A,v,P,L}$; for any $t \ge 1$, if (f_t, g_t) is the response of the procedure at time t for query x_t , then $f_t = F_{A,v,P,L}(x_t)$ and $g_t = \partial F_{A,v,P,L}(x_t)$.

Proof. Consider a given iteration t. We aim to show that $(f_t, \mathbf{g}_t) = (F_{A,v,P,L}(\mathbf{x}_t), \partial F_{A,v,P,L}(\mathbf{x}_t))$. By construction, if $t \ge d^2$, the result is immediate. Now, suppose $t \le d^2$. We first consider the case when \mathbf{x}_t is noninformative (1). By definition, $F_A(\mathbf{x}_t) > \eta$. Because for any $(p,l) \le_{lex} (P,L)$, one has $|v_{p,l}^\top \mathbf{x}_t| \le ||v_{p,l}|| ||x_t|| \le 1$, we have

$$F_{A,v,P,L}(\mathbf{x}_t) = \max \left\{ F_A(\mathbf{x}_t), \eta \left(\max_{(p,l) \leq lex(P,L)} \mathbf{v}_{p,l}^\top \mathbf{x} - p \gamma_1 - l \gamma_2 \right) \right\} = F_A(\mathbf{x}_t).$$

As a result, the response of the procedure for x_t is consistent with $F_{A,v,P,L}$, and the returned subgradient is $\tilde{g}_A(x_t) = \partial F_{A,v,P,L}(x_t)$. Therefore, it suffices to focus on informative queries (2). We will denote by $t_{p,l}$ the index of the iteration when $v_{p,l}$ has been defined for $(p,l) \leq_{lex}(P,L)$. Consider a specific couple $(p,l) \leq_{lex}(P,L)$, and let r denote the number of constructed indices on or before $t_{p,l}$. Let $b_{p,1},\ldots,b_{p,r}$ be the corresponding vectors resulting from the Gram–Schmidt procedure on $x_{i_{p,1}},\ldots,x_{i_{p,r}}$. Then, conditionally on the history until time $t_{p,l}$, the vector $v_{p,l}$ was defined as $v_{p,l} = \phi_\delta(y_{p,l})$, where $y_{p,l}$ is sampled as $\sim \mathcal{U}(S^{d-1} \cap \{z \in \mathbb{R}^d : |b_{p,r'}^\top z| \leq d^{-3}, \forall r' \leq r\})$. As a result, from Lemma A.1, for any $t \leq t_{p,l}$, we have

$$\mathbb{P}\left(|\boldsymbol{x}_t^{\top}\boldsymbol{v}_{p,l}| \geq 3\sqrt{\frac{2\log d}{d}} + \frac{2}{d^2}\right) \leq \frac{6\sqrt{2\log d}}{d^6}.$$

We then define the following event

$$\mathcal{E} = \bigcap_{(p,l) \leq_{lex}(P,L)} \bigcap_{t \leq t_{p,l}} \left\{ |\mathbf{x}_t^\top \mathbf{v}_{p,l}| < 3\sqrt{\frac{2\log d}{d}} + \frac{2}{d^2} \right\},\,$$

which by the union bound, has probability $\mathbb{P}(\mathcal{E}) \ge 1 - 3\sqrt{2\log d}/d^2$. We are now ready to show that the construction procedure is consistent with optimizing $F_{A,v,P,L}$ on the event \mathcal{E} . As seen before, we can suppose that x_t is informative (2). Using the same notations as before, because \mathcal{E} is met, for any $p < p' \le P$ and $l' \le l_{p'}$, we have for $d \ge 2$,

$$\boldsymbol{v}_{p',l'}^{\top}\boldsymbol{x}_t - p'\boldsymbol{\gamma}_1 - l'\boldsymbol{\gamma}_2 < 3\sqrt{\frac{2\log d}{d}} + \frac{1}{d} - p\boldsymbol{\gamma}_1 - \boldsymbol{\gamma}_1 \leq -p\boldsymbol{\gamma}_1 - \frac{\boldsymbol{\gamma}_1}{2} \leq -p\boldsymbol{\gamma}_1 - d\boldsymbol{\gamma}_2 - \frac{\boldsymbol{\gamma}_2}{2},$$

where we used $3\sqrt{2} + 1 \le 6$ and $2d\gamma_2 \le \gamma_1/2$. As a result, we obtain that

$$\max_{(p',l')\leq_{lex}(P,L),\,p'>p}v_{p',l'}^{\intercal}x_t-p'\gamma_1-l'\gamma_2<-p\gamma_1-l\gamma_2-\frac{\gamma_2}{2}.$$

Next, we consider the case of vectors $v_{p,l'}$, where $l \le l' \le l_p$ and $t_{p,l'} \ge t$ (this also includes the case when we defined $v_{p,l}$ at time $t = t_{p,l}$). We write \tilde{l} for the smallest such index l. As a remark, $\tilde{l} \in \{l, l+1\}$. Note that if such indices exist, this means that before starting iteration t, the procedure has not yet reached r = k. There are two cases. If x_t was exploratory, we have $t = i_{p,r}$; hence, $\|P_{Span(b_n,r,r'\le r)^\top}(x_t)\| = 0$. If x_t is not exploratory, either

$$||P_{Span(b_{p,r'},r'\leq r)^{\top}}(x_t)|| < \frac{\gamma_2}{4}||x_t|| \le \frac{\gamma_2}{4},$$
 (5)

or we have $F_{A,v,p,l}(x_t) > -\eta \gamma_1/2$. We start with the last scenario when $F_{A,v,p,l}(x_t) > -\eta \gamma_1/2$. Then, on \mathcal{E} , one has

$$\max_{(p,l) <_{lex}(p',l') \leq_{lex}(P,L)} v_{p',l'}^{\top} x_t - p' \gamma_1 - l' \gamma_2 \leq -\gamma_1 + 3\sqrt{\frac{2\log d}{d}} + \frac{1}{d} \leq -\frac{\gamma_1}{2}.$$

As a result, this shows that $F_{A,v,P,L}(x_t) = F_{A,v,p,l}(x_t)$. Hence, using a first-order oracle from $F_{A,v,l,p}$ at x_t is already consistent with $F_{A,v,P,L}$. Thus, for whichever case (case (2a), case (2b), or case (2c)) is performed, because these can only increase the knowledge on v, the response given by the construction procedure is consistent with minimizing $F_{A,v}$.

It remains to treat the first two scenarios in which we always have Equation (5). In particular, when writing $x_t = \alpha_1 \boldsymbol{b}_{p,1} + \dots + \alpha_r \boldsymbol{b}_{p,r} + \tilde{\boldsymbol{x}}_t$ where $\tilde{\boldsymbol{x}}_t = P_{Span(\boldsymbol{b}_{p,r},r' \leq r)^{\perp}}(\boldsymbol{x}_t)$, we have $\|\tilde{\boldsymbol{x}}_t\| < \gamma_2/4$. As a result, for $\tilde{l} \leq l' \leq l_p$, one has for

$$\begin{split} |\boldsymbol{v}_{p,l'}^{\top}\boldsymbol{x}_{t}| &\leq |\boldsymbol{y}_{p,l'}^{\top}\boldsymbol{x}_{t}| + \delta \leq |\alpha_{1}||\boldsymbol{y}_{p,l'}^{\top}\boldsymbol{b}_{p,1}| + \dots + |\alpha_{r}||\boldsymbol{y}_{p,l'}^{\top}\boldsymbol{b}_{p,r}| + ||\tilde{\boldsymbol{x}}_{t}|| + \delta \\ &< ||\boldsymbol{\alpha}||_{1}\frac{1}{d^{3}} + \frac{\gamma_{2}}{4} + \delta \\ &\leq \frac{\gamma_{2}}{4} + \frac{1}{d^{2}\sqrt{d}} + \frac{1}{d^{3}} \leq \frac{\gamma_{2}}{2}, \end{split}$$

where in the last inequality, we used $d \ge 3$. As a result, provided that \hat{l} exists, this shows that

$$\max_{\tilde{l} \le l' \le l_n} v_{p,l'}^{\top} x_t - p \gamma_1 - l' \gamma_2 = v_{p,\tilde{l}}^{\top} x_t - p \gamma_1 - \tilde{l} \gamma_2 < -p \gamma_1 - \tilde{l} \gamma_2 + \frac{\gamma_2}{2}.$$
 (6)

On the other hand, if $t = i_{p+1,1}$, the same reasoning works for t viewing it as in period p + 1, which shows for this case that

$$\max_{l' \le l_{p+1}} v_{p+1,l'}^{\top} x_t - (p+1)\gamma_1 - l' \gamma_2 = v_{p+1,1}^{\top} x_t - (p+1)\gamma_1 - \gamma_2 < -(p+1)\gamma_1 - \frac{\gamma_2}{2}. \tag{7}$$

As a conclusion of these estimates, we showed that on \mathcal{E} , we have

$$F_{A,v,P,L}(\mathbf{x}_t) = \max\{F_{A,v,p,l}(\mathbf{x}_t), \eta(\mathbf{v}_{v',l'}^{\top}\mathbf{x}_t - p'\gamma_1 - l'\gamma_2)\} := \tilde{F}_{A,v,t}(\mathbf{x}_t),$$

where (p',l') is the very next vector that is defined after starting iteration t (potentially, it has $t_{p',l'}=t$ if we defined a vector at this time). It now suffices to check that the value and the vector returned by the procedure are consistent with the right-hand side. By construction, if we constructed $v_{p',l'}$ at step t (case (2b) or case (2c)), then the procedure directly uses a first-order oracle for $\tilde{F}_{A,v,t}$. Further, by construction of the subgradients because they break ties lexicographically in (p,l), the returned subgradient is exactly $\partial F_{A,v,P,L}(x_t)$. It remains to check that this is the case when no vector $v_{p',l'}$ is defined at step t: case (2a). This corresponds to the case when $F_{A,v,p,l}(x_t) \ge \eta(-p\gamma_1 - l\gamma_2 - \gamma/2)$. Now, in this case, the upper-bound estimates from Equations (6) and (7) imply that

$$\boldsymbol{v}_{p',l'}^{\top}\boldsymbol{x}_t - p'\boldsymbol{\gamma}_1 - l'\boldsymbol{\gamma}_2 < -p\boldsymbol{\gamma}_1 - l\boldsymbol{\gamma}_2 - \boldsymbol{\gamma}/2,$$

and as a result, $F_{A,v,P,L}(x_t) = F_{A,v,p,l}(x_t)$. Therefore, using a first-order oracle of $F_{A,v,P,L}$ at x_t is valid, and the break of ties of the subgradient of $\tilde{F}_{A,v,t}$ is the same as the break of ties of $\partial F_{A,v,P,L}(x_t)$. This ends the proof that on \mathcal{E} , the procedure gives responses consistent with an optimization oracle for $F_{A,v,P,L}$ with subgradient function $\partial F_{A,v,P,L}$. Because $\mathbb{P}(\mathcal{E}) \geq 1 - C\sqrt{\log d}/d^2$ for some constant C > 0, this ends the proof of the proposition. \Box

Last, we provide an upper bound on the optimal value of $F_{A,v,P,L}$.

Proposition 2. Let $A \sim \mathcal{U}(\{\pm 1\}^{n \times d})$ and $v_0 \sim \mathcal{U}(\mathcal{D}_{\delta})$. For any algorithm alg for convex optimization, let v be the resulting set of vectors constructed by the randomized procedure. With probability at least $1 - C\sqrt{\log d}/d$ over the randomness of A, v_0 and v, we have

$$\min_{x \in B_d(0,1)} F_{A,v}(x) \le -\frac{\eta}{40\sqrt{(kp_{max} + 1)\log d}}$$

for some universal constant C > 0.

Proof. For simplicity, let us enumerate all of the constructed vectors $v_1, \ldots, v_{l_{max}}$ by order of construction. Hence, $l_{max} \leq p_{max}k$. We use the same enumeration for $y_1, \ldots, y_{l_{max}}$. Now, let $C_d = \sqrt{40(l_{max}+1)\log d}$, and consider the following vector:

$$\bar{\boldsymbol{x}} = -\frac{1}{C_d} \sum_{l=0}^{l_{max}} P_{Span(\boldsymbol{a}_l, i \leq n)^{\perp}}(\boldsymbol{v}_l).$$

In particular, note that we included v_0 in the sum. For convenience, we write $P_{A^{\perp}}$ instead of $P_{Span(a_l,i\leq n)^{\perp}}$. Also, for convenience, let us define $z_l = \sum_{l'\leq l} P_{A^{\perp}}(v_l)$. Fix an index $1\leq l\leq l_{max}$. Then, by Lemma A.1, with $t_0:=\sqrt{6\log d/d}+2d^{-2}$, we have

$$\begin{split} \mathbb{P}(|P_{A^{\perp}}(v_{l+1})^{\top}z_{l}| > t_{0}||z_{l}||) &= \mathbb{P}(|v_{l+1}^{\top}P_{A^{\perp}}(z_{l})| > t_{0}||z_{l}||) \\ &\leq \mathbb{P}(|v_{l+1}^{\top}P_{A^{\perp}}(z_{l})| > t_{0}||P_{A^{\perp}}(z_{l})||) \\ &\leq \frac{2\sqrt{6\log d}}{d^{2}}. \end{split}$$

Similarly, we have that

$$\mathbb{P}(|v_{l+1}^{\top} z_l| > t_0 ||z_l||) \le \frac{2\sqrt{6 \log d}}{d^2}.$$

Now, consider the event $\mathcal{E} = \bigcap_{l \leq l_{max}} \{ |v_l^\top z_{l-1}|, |P_{A^\perp}(v_l)^\top z_{l-1}| \leq t_0 ||z_l|| \}$, which because $l_{max} \leq d$, by the union bound has probability at least $1 - 4\sqrt{6\log d}/d$. Then, on \mathcal{E} , for any $l < l_{max}$,

$$||z_{l+1}||^2 \le ||z_l||^2 + ||P_{A^{\perp}}(v_{l+1})||^2 + 2|P_{A^{\perp}}(v_{l+1})^{\top}z_l| \le ||z_l||^2 + 1 + 2t_0||z_l||.$$

We now prove by induction that $||z_l||^2 \le 40 \log d \cdot (l+1)$, which is clearly true for z_0 because $||z_0|| = ||P_{A^{\perp}}(v_0)|| \le ||v_0|| \le 1$. Suppose this is true for $l < l_{max}$. Then, using the above equation and the fact that $t_0 \le 3\sqrt{\log d/d}$ for $d \ge 4$,

$$||z_{l+1}||^2 \le 40 \log d \cdot (l+1) + 1 + 6\sqrt{40} \log d \sqrt{\frac{l+1}{d}} \le 40 \log d \cdot (l+2),$$

where we used $l_{max} + 1 \le d$, which completes the induction. In particular, on \mathcal{E} , we have that $||\bar{x}|| \le 1$. Now, observe that by construction, $\bar{x} \in Span(a_i, i \le n)^{\perp}$ so that $||A\bar{x}||_{\infty} = 0$. Next, for any $0 \le l \le l_{max}$, we have

$$\boldsymbol{v}_{l}^{\top} \bar{\boldsymbol{x}} = -\frac{\boldsymbol{v}_{l}^{\top} \boldsymbol{z}_{l_{max}}}{C_{d}} = -\frac{1}{C_{d}} \left(\|P_{A^{\perp}}(\boldsymbol{v}_{l})\|^{2} + \boldsymbol{v}_{l}^{\top} \boldsymbol{z}_{l-1} + \sum_{l < l' \leq l_{max}} \boldsymbol{v}_{l}^{\top} P_{A^{\perp}}(\boldsymbol{v}_{l'}) \right).$$

We will give estimates on each term of the above equation. First, if the indices $i_{p,1}, \ldots, i_{p,r}$ were defined before defining v_l , we denote $\tilde{y} = P_{Span(x_{i_{p,r}}, r' \leq r)^{\perp}}(y_l)$, the component of y_l that is perpendicular to the explored space at that time. Then, we can write $y_l = \alpha_1^l b_{p,1} + \cdots + \alpha_r^l b_{p,1} + \tilde{y}_l$ and note that

$$\|\tilde{y}_l\| = \sqrt{\|y_l\| - (\alpha_1^l)^2 - \dots - (\alpha_r^l)^2} \ge \sqrt{1 - \frac{k}{d^6}} \ge 1 - \frac{1}{d^5}$$

Then, we have

$$\begin{split} \|P_{A^{\perp}}(\boldsymbol{v}_l)\| &\geq \|P_{A^{\perp}}(\boldsymbol{y}_l)\| - \delta \\ &\geq \|P_{Span(a_i, i \leq n, \ \boldsymbol{b}_{p,r'}, r \leq r')^{\perp}}(\boldsymbol{y}_l)\| - \delta \\ &= \|P_{Span(a_i, i \leq n, \ \boldsymbol{b}_{p,r'}, r \leq r')^{\perp}}(\tilde{\boldsymbol{y}}_l)\| - \delta \\ &\geq \left\|P_{Span(a_i, i \leq n, \ \boldsymbol{b}_{p,r'}, r' \leq r)^{\perp}}\left(\frac{\tilde{\boldsymbol{y}}_l}{\|\tilde{\boldsymbol{y}}_l\|}\right)\right\| - \frac{1}{d^5} - \delta. \end{split}$$

As a result, because $\delta = d^{-3}$, this shows that

$$||P_{A^{\perp}}(\boldsymbol{v}_l)||^2 \ge \left||P_{Span(a_i, i \le n, \, \boldsymbol{b}_{p,r'}, r' \le r)^{\perp}}\left(\frac{\tilde{\boldsymbol{y}}_l}{||\tilde{\boldsymbol{y}}_l||}\right)\right||^2 - 2\delta.$$

Now, observe that $dim(Span(a_i, i \le n, b_{p,r'}, r' \le r)^{\perp}) \ge d - n - k$, whereas $\tilde{y}_l / ||\tilde{y}_l||$ is a uniformly random unit vector in $Span(b_{p,r'}, r \le r')^{\perp}$. Therefore, using Proposition A.1, we obtain for t < 1,

$$\begin{split} & \mathbb{P}\bigg(\|P_{A^{\perp}}(v_l)\|^2 + 2\delta - \frac{d-n-k}{d} \leq -t\bigg) \\ & \leq \mathbb{P}\bigg(\bigg\|P_{Span(a_i,i\leq n,\; \boldsymbol{b}_{p,r'},\, r'\leq r)^{\perp}}\bigg(\frac{\tilde{y}_l}{\|\tilde{y}_l\|}\bigg)\bigg\|^2 - \frac{d-n-k}{d} \leq -t\bigg) \\ & \leq e^{-(d-k)t^2}. \end{split}$$

As a result, because $d - n - k \ge d/2$, we obtain

$$\mathbb{P}\left(\left\|P_{A^{\perp}}(v_l)\right\|^2 \leq \frac{1}{2} - 2\sqrt{\frac{\log d}{d}} - 2\delta\right) \leq \frac{1}{d^2}.$$

Now, define $\mathcal{F} = \bigcap_{l \le l_{max}} \{ \|P_{A^{\perp}}(v_l)\|^2 \ge 1/2 - 2\sqrt{\log d/d} - 2\delta \}$, which because $l_{max} + 1 \le d$ and by the union bound

has probability at least $\mathbb{P}(\mathcal{F}) \geq 1 - 1/d$. Next, we turn to the last term. For any $0 \leq l < l_{max}$, we now focus on the sequence $(\sum_{l'=l+1}^{l+u} v_l^{\top} P_{A^{\top}}(y_{l'}))_{1 \leq u \leq l_{max}-l}$ and first note that this is a martingale. These increments are symmetric (because $y_{l'}$ is symmetric) even conditionally on A and $v_l, y_l, \ldots, y_{l'-1}$. Next, let $t_1 = 2\sqrt{3\log d/d} + 2d^{-2}$. Note that for $d \geq 4$, we have $t_1 \leq 4\sqrt{\log d/d}$. Further, by Lemma A.1,

$$\mathbb{P}(|v_l^{\top} P_{A^{\top}}(y_{l'})| > t_1) = \mathbb{P}(|P_{A^{\top}}(v_l)^{\top} y_{l'}| > t_1) \leq \frac{4\sqrt{3 \log d}}{d^4}$$

where we used the fact that $P_{A^{\perp}}$ is a projection. Let $\mathcal{G}_l = \bigcap_{l < l' \le l_{max}} \{|v_l^{\top} P_{A^{\top}}(v_{l'})| \le t_1\}$, which by the union bound has probability $\mathbb{P}(\mathcal{G}_l) \ge 1 - 4\sqrt{3\log d}/d^3$. Next, we define $I_{l,u} = (v_l^{\top} P_{A^{\top}}(y_{l+u}) \wedge t_1) \vee (-t_1)$, the increments capped at absolute value t_1 . Because $v_l^{\top} P_{A^{\top}}(y_{l+u})$ is symmetric, so is $I_{l,u}$. As a result, these are bounded increments of a martingale, to which we can apply the Azuma–Hoeffding inequality:

$$\mathbb{P}\left(\left|\sum_{u=1}^{l_{max}-l}I_{l,u}\right| \leq 2t_1\sqrt{(l_{max}-l)\log d}\right) \geq 1-\frac{2}{d^2}.$$

We denote by \mathcal{H}_l this event. Now, observe that on \mathcal{G}_l , the increments $I_{l,u}$ and $v_l^{\mathsf{T}} P_{A^{\mathsf{T}}}(y_{l+u})$ coincide for all $1 \leq u \leq l_{max} - l$. As a result, on $\mathcal{G}_l \cap \mathcal{H}_l$, we obtain

$$\left| \sum_{l < l' \le l_{max}} \boldsymbol{v}_l^{\top} P_{A^{\perp}}(\boldsymbol{v}_{l'}) \right| \le \left| \sum_{l < l' \le l_{max}} \boldsymbol{v}_l^{\top} P_{A^{\perp}}(\boldsymbol{y}_{l'}) \right| + (l_{max} - 1)\delta$$

$$\le \left| \sum_{u=1}^{l_{max} - l} l_{l,u} \right| + (d - 2)\delta$$

$$\le 2t_1 \sqrt{l_{max} \log d} + (d - 2)\delta.$$

Then, on the event $\mathcal{E} \cap \mathcal{F} \cap \cap_{l \leq l_{max}} \mathcal{G}_l \cap \mathcal{H}_l$, for any $1 \leq l \leq l_{max}$, one has

$$\begin{split} v_l^\top z_{l_{max}} &\geq \frac{1}{2} - 2\sqrt{\frac{\log d}{d}} - t_0 ||z_l|| - 2t_1 \sqrt{l_{max} \log d} - \frac{1}{d^2} \\ &\geq \frac{1}{2} - 2\sqrt{\frac{\log d}{d}} - 3\log d\sqrt{40\frac{l_{max} + 1}{d}} - 8\log d\sqrt{\frac{l_{max}}{d}} - \frac{1}{d^2} \\ &\geq \frac{1}{2} - 30\log d\sqrt{\frac{l_{max} + 1}{d}} \\ &\geq \frac{1}{6}, \end{split}$$

where in the last inequalities, we used the fact that $l_{max} \le kp_{max} \le c_{d,1}d - 1$, where $c_{d,1} = 1/(90^2 \log^2 d)$ as per Equation (4). As a result, we obtain that on $\mathcal{E} \cap \mathcal{F} \cap \cap_{l \le l_{max}} \mathcal{G}_l \cap \mathcal{H}_l$, which has probability at most $1 - C\sqrt{\log d}/d$ for some constant C > 0,

$$\max_{p \le p_{max}, l \le k} v_{p, l}^{\top} \bar{x} \le -\frac{1}{6C_d} \le -\frac{1}{40\sqrt{(kp_{max} + 1)\log d}}.$$

Because $||A\bar{x}||_{\infty} = 0$ and $\eta \ge \eta/(40\sqrt{(kp_{max} + 1)\log d})$, this shows that

$$F_{A,v}(\bar{x}) \le -\frac{\eta}{40\sqrt{(kp_{max}+1)\log d}}.$$

This ends the proof of the proposition. \Box

3.3. Reduction from Convex Optimization to the Optimization Procedure

According to Proposition 1, with probability at least $1-C\sqrt{\log d}/d^2$, the procedure returns responses that are consistent with a first-order oracle of the function $F_{A,v,P,L}$, where $v_{P,L}$ is the last vector to have been defined. Now, observe that for any constructed vectors v, the function $F_{A,v,P,L}$ is \sqrt{d} -Lipschitz. As a result, if there exists an algorithm for convex optimization that guarantees ϵ accuracy for 1-Lipschitz functions, by rescaling, there exists an algorithm alg that is successful for the optimization procedure with probability $1-C\sqrt{\log d}/d^2$ and $\epsilon\sqrt{d}$ accuracy. In the next proposition, we show that to be successful, such an algorithm needs to properly define the complete function $F_{A,v}$ (i.e., to complete all periods until p_{max}).

Proposition 3. Let alg be a successful algorithm for the optimization procedure with probability $q \in [0,1]$ and precision $\eta/(2\sqrt{d})$. Suppose that alg performs at most d^2 queries during the optimization procedure. Then, when running alg with the responses of the optimization procedure, alg succeeds and ends the period p_{max} with probability at least $q - C\sqrt{\log d}/d$ for some universal constant C > 0.

Proof. Let $x^*(alg) = x_T$ denote the final answer of *alg* when run with the optimization procedure. By hypothesis, we have $T \le d^2$. As before, let $P \le p_{max}$ and $L \le k$ be the indices such that the last vector constructed by the optimization procedure is $v_{P,L}$. Let \mathcal{E} be the event when *alg* run on the optimization procedure does not end period p_{max} . We focus on \mathcal{E} and consider two cases.

First, suppose that $T > t_{P,L}$ (i.e., the last vector was not constructed at time T). As a result, either this means that x_T corresponds to a noninformative query—case (1)—in which case $F_{A,v,P,L}(x_T) \ge F_A(x_T) \ge \eta$, or this means that $F_{A,v,P,L}(x_t) \ge \eta(-P\gamma_1 - L\gamma_2 - \gamma/2)$ —case (2a).

Second, we now suppose that $T = t_{P,L}$ (i.e., the last vector was constructed at time T). Then, by construction of $v_{P,L}$ and $y_{P,L}$, we have indices $i_{P,1}, \ldots, i_{P,r} \leq T$ such that with the Gram–Schmidt decomposition $b_{P,1}, \ldots, b_{P,r}$ of $x_{i_{P,1}}, \ldots, x_{i_{P,r}}$, we have $|b_{p,r}^{\top}, y_{P,L}| \leq d^{-3}$ for all $r' \leq r$. In particular, writing $x_T = \alpha_1 b_{P,1} + \cdots + \alpha_r b_{P,r} + \tilde{x}_T$, where $\tilde{x}_T \in Span(x_{i_{P,r}}, r' \leq r)^{\perp}$, either we have $i_{P,r} = T$, in which case $\tilde{x}_T = \mathbf{0}$, or x_T was not exploratory, in which case we directly have $F_{A,v,P,L}(x_T) \geq F_{A,v,P,L-1}(x_T) > -\eta \gamma_1/2$, or we have $||\tilde{x}_T|| < ||x_T|| \gamma_2/4 \leq \gamma_2/4$. For all remaining cases to consider, we obtain

$$|v_{P,L}^{\top}x_{T}| \leq |y_{P,L}^{\top}x_{T}| + \delta \leq \frac{\|\alpha\|_{1}}{d^{3}} + \|\tilde{x}_{T}\| + \delta \leq \frac{1}{d^{3}} + \frac{1}{d^{2}\sqrt{d}} + \frac{\gamma_{2}}{4} < \frac{\gamma_{2}}{2}.$$

In the last inequality, we used $d \ge 4$. This shows that $F_{A,v,P,L}(\mathbf{x}_T) \ge \eta(-P\gamma_1 - L\gamma_2 - \gamma_2/2)$. As a result, in all cases, this shows that $F_{A,v,P,L}(\mathbf{x}^*(alg)) \ge \eta(-P\gamma_1 - L\gamma_2 - \gamma_2/2) \ge -\eta(p_{max} + 1)\gamma_1$. Now, define the event

$$\mathcal{F} = \left\{ \min_{x \in B_d(0,1)} F_{A,v}(x) \le -\frac{\eta}{40\sqrt{(kp_{max} + 1)\log d}} \right\}.$$

By Proposition 2, we have $\mathcal{P}(\mathcal{F}) \ge 1 - C\sqrt{\log d}/d$. Now, from Equation (4),

$$(p_{max} + 1)^{3/2} \le \frac{1}{60\gamma_1 \sqrt{k \log d}}.$$

Thus,

$$(p_{max} + 1)\gamma_1 \le \frac{1}{60\sqrt{k(p_{max} + 1)\log d}} \le \frac{1}{60\sqrt{(kp_{max} + 1)\log d}}.$$

Then, because $F_{A,v,P,L} \leq F_{A,v}$, this shows that on $\mathcal{E} \cap \mathcal{F}$,

$$\begin{split} F_{A,v,P,L}(x^{\star}(alg)) &\geq -\eta(p_{max}+1)\gamma_{1} \geq \min_{x \in B_{d}(0,1)} F_{A,v}(x) + \frac{\eta}{120\sqrt{(kp_{max}+1)\log d}} \\ &> \min_{x \in B_{d}(0,1)} F_{A,v,P,L}(x) + \frac{\eta}{2\sqrt{d}}, \end{split}$$

where in the last inequality, we used $kp_{max} \leq c_{d,1}d - 1$. As a result, let \mathcal{G} be the event when alg succeeds for precision $\epsilon = \eta/(2\sqrt{d})$. By hypothesis, $\mathcal{P}(\mathcal{G}) \geq q$. Now, from the above equations, one has $\mathcal{E} \cap \mathcal{F} \cap \mathcal{G} = \emptyset$. Therefore, $\mathbb{P}(\mathcal{G} \cap \mathcal{E}^c) \geq \mathcal{P}(\mathcal{G}) - \mathbb{P}(\mathcal{G} \cap \mathcal{E} \cap \mathcal{F}) - \mathbb{P}(\mathcal{F}^c) \geq q - C\sqrt{\log d}/d$. This ends the proof of the proposition. \square

3.4. Reduction from an Orthogonal Vector Game with Hints to the Optimization Procedure

We are now ready to introduce an orthogonal vector game, where the main difference with the game introduced in Marsden et al. [22] is that the player can provide additional hints. The game is formally described in Game 1.

Let us first describe the game from Marsden et al. [22]. At the start of the game, an oracle samples a random matrix $A \sim \mathcal{U}(\{\pm 1\}^{n \times d})$, where the number of rows n is fixed (say, $n = \lfloor d/4 \rfloor$). The final goal of the player is to output k vectors y_1, \ldots, y_k , where k is a parameter of the game, that are approximately orthogonal to A and also robustly linearly independent. Precisely, for some parameters α , β of the game, these should satisfy

$$\|A\boldsymbol{y}_i\|_{\infty} \leq \alpha \quad \text{and} \ \|P_{Span(\boldsymbol{y}_1,\dots,\boldsymbol{y}_{i-1})^\perp}(\boldsymbol{y}_i)\|_2 \geq \beta, \quad i \in [k].$$

To find these vectors, the player can only use an *M*-bit message Message that was previously constructed with the knowledge of *A* and query *m* rows of *A*. Marsden et al. [22] showed that to solve this game, a player either needs a large memory for Message or needs to query most of the rows of *A*.

In Game 1, the player has further access to *hints*, which are vectors v_1, \ldots, v_l that have also been constructed previously using the knowledge of A. In this game, when making queries in line 9–12, the player first submits a vector $z \in \mathbb{R}^d$ and then receives a response g(z), where g is a function that takes as values either the rows of A or the hints v_1, \ldots, v_d and that was also previously specified by the player. For technical reasons, we also allow the response g to return an additional number in $[d^2]$ (see line 8 of Game 1; in any case, this number will carry little information about A). The game, therefore, has two phases. First, knowing A, the player chooses Message, the hints v_1, \ldots, v_d , and a response function g (lines 2–8 of Game 1). Second, the player aims to output solution vectors y_1, \ldots, y_k using only Message and m queries to g.

Note that if the hints $v_1, ..., v_d$ could be specified by the player without constraints, the player could directly choose as hints some solution vectors $v_1, ..., v_k$ that are orthogonal to A and robustly independent. Then, the player would easily win in the second phase with at most m = k queries to g. Instead, the hints $v_1, ..., v_d$ are constructed in a very specific way that exactly mimics the construction of the vectors $v_{p,l}$ in Procedure 1. To construct each vector v_l in the first phase, the player submits at most k vectors $x_{l,1}, ..., x_{l,r_l}$ (these mimic exploratory queries). The vector v_l is then obtained after discretizing a random vector approximately orthogonal to $x_{l,1}, ..., x_{l,r_l}$ as in lines 12 and 13 of Procedure 1. The goal in the next sections will be to show that even with the additional information from the hints, the game is still hard to win for the player.

Game 1 (Orthogonal Vector Game with Hints)

Input: d, k, m, M, α , β

- 1 Oracle: Set $n \leftarrow \lfloor d/4 \rfloor$, sample $A \sim \mathcal{U}(\{\pm 1\}^{n \times d})$
- 2 Player: Observe A
- 3 for $l \in [d]$, do
- 4 Player: Based on A and any previous queries and responses, submit at most k vectors $x_{l,1}, \dots, x_{l,r_l}$
- Oracle: Perform the Gram–Schmidt decomposition $b_{l,1},\ldots,b_{l,r_l}$ of $x_{l,1},\ldots,x_{l,r_l}$. Then, sample a vector $y_l \in S^{d-1}$ according to a uniform distribution $\mathcal{U}(S^{d-1} \cap \{z \in \mathbb{R}^d : \forall r \leq r_l, |b_{l,r}^{\top}z| \leq d^{-3}\})$. As response to the query, return $v_l = \phi_{\delta}(y_l)$ to the player.
- 6 end
- 7 Player: Based on A, all previous queries and responses, store an M-bit message Message
- 8 *Player*: Based on *A*, all previous queries and responses, submit a function $g: B_d(0,1) \to (\{a_j, j \le n\} \cup \{v_l, l \le d\}) \times [d^2]$ to the oracle
- 9 for $i \in [m]$, do
- 10 *Player*: Based on Message and any previous queries $x_1,...,x_{i-1}$ and responses $g_1,...,g_{i-1}$ from this loop phase, submit a query $z_i \in \mathbb{R}^d$
- 11 *Oracle*: As the response to query z_i , return $g_i = g(z_i)$
- 12 end
- 13 *Player*: Based on all queries and responses from this phase $\{z_i, g_i, i \in [m]\}$ and on Message, return some vectors y_1, \ldots, y_k to the oracle
- 14 The player wins if the returned vectors have unit norm and satisfy for all $i \in [k]$
 - $1 \|Ay_i\|_{\infty} \leq \alpha$
 - 2 $||P_{Span(y_1,...,y_{i-1})^{\perp}}(y_i)||_2 \ge \beta$

We first prove that solving the optimization procedure implies solving the orthogonal vector game with hints.

Proposition 4. Let $m \le d$. Suppose that there is an M-bit algorithm that is successful for the optimization procedure with probability q for accuracy $\epsilon = \eta/(2\sqrt{d})$ and uses at most mp_{max} queries. Then, there is an algorithm for Game 1 for parameters $(d,k,m,M,\alpha=2\eta/\gamma_1,\beta=\gamma_2/4)$, for which the player wins with probability at least $q-C\sqrt{\log d}/d$ for some universal constant C>0.

Proof. Let *alg* be an *M*-bit algorithm solving the feasibility problem with mp_{max} queries with probability at least q. We now describe the strategy for Game 1.

In the first part of the strategy, the player observes A. First, submit an empty query to the oracle to obtain a vector v_0 , which as a result, is uniformly distributed among \mathcal{D}_{δ} . We then proceed to simulate the optimization procedure for alg using parameters A and v_0 (lines 3–6 of Game 1). Precisely, whenever a new vector $v_{p,l}$ needs to be defined according to the optimization procedure, the player submits the corresponding vectors $x_{i_{p,1}},\ldots,x_{i_{p,r}}$ to the oracle and receives in return a vector that defines $v_{p,l}$. In this manner, the player simulates exactly the optimization procedure. In all cases, the number of queries in this first phase is at most $1+kp_{max}\leq d$. For the remaining queries to perform, the player can query whichever vectors; these will not be used in the rest of the strategy. If the simulation did not end period p_{max} , the complete procedure fails. We now describe the rest of the procedure when period p_{max} was ended. During the simulation, the algorithm records the time $i_{p,1}$ when period p_{max} started for all $p \leq p_{max} + 1$. Recall that for $p_{max} + 1$, we only define $i_{p_{max}+1,1}$; this is the time that ends period p_{max} . Now, by hypothesis, $i_{p_{max}+1,1} \leq mp_{max}$. As a result, there must be a period $p \leq p_{max}$ that uses at most m queries: $i_{p+1,1} - i_{p,1} \leq m$. We define the memory Message to be the memory of alg just before starting iteration $i_{p,1}$ at the beginning of period p (line 7 of Game 1). Next, because the period p_{max} was ended, the vectors $v_{p,l}$ for $p \leq p_{max}, l \leq l_p$ were all defined. The player can, therefore, submit the function $g_{A,v}$ to the oracle (line 8 of Game 1) as follows:

$$g_{A,v}: x \longmapsto \begin{cases} (g_A(x), 1) & \text{if } F_{A,v}(x) = ||Ax||_{\infty} - \eta, \\ (v_0, 2) & \text{otherwise and if } F_{A,v}(x) = \eta v_0^{\mathsf{T}} x, \\ (v_{p,l}, 2 + (p-1)k + l) & \text{otherwise and if} \\ (p,l) = \underset{(p',l') \leq_{lex}(p_{max}, l_{p_{max}})}{\arg \max} v_{p',l'}^{\mathsf{T}} x - p \gamma_1 - l \gamma_2. \end{cases}$$

$$(8)$$

Intuitively, the first component of $g_{A,v}$ gives the subgradient $\partial F_{A,v}$ to the following two exceptions; we always return a_i instead of $\pm a_i$, and we return v_0 ($v_{p,l}$, respectively) instead of ηv_0 ($\eta v_{p,l}$, respectively). The second term of $g_{A,v}$ has values in $[2 + p_{max}k]$. Hence, because $2 + p_{max}k \le d^2$, the function $g_{A,v}$ takes values in $(\{a_j, j \le n\} \cup \{v_l, l \le d\}) \times [d^2]$.

The strategy then proceeds to play the orthogonal vector game in a second part (lines 9–12 of Game 1) and uses the responses of the oracle to simulate the run of alg for the optimization procedure in period p. To do so, we set the memory state of the algorithm alg to be Message. Then, for the next m iterations, we proceed as follows. At iteration i of the process, we run alg with its current state to obtain a new query z_i , which is then submitted to the oracle of the orthogonal vector game to get a response (g_i, s_i) . We then use this response to simulate the response that was given by the optimization procedure in the first phase, computing (v_i, \tilde{g}_i) as follows:

$$(v_i, \tilde{\mathbf{g}}_i) = \begin{cases} (|\mathbf{g}_i^{\mathsf{T}} \mathbf{z}_i| - \eta, sign(\mathbf{g}_i^{\mathsf{T}} \mathbf{z}_i) \mathbf{g}_i) & s_i = 1, \\ (\eta \mathbf{g}_i^{\mathsf{T}} \mathbf{z}_i, \eta \mathbf{g}_i) & s_i = 2, \\ (\eta (\mathbf{g}_i^{\mathsf{T}} \mathbf{z}_i - p \gamma_1 - l \gamma_2), \eta \mathbf{g}_i) & s_i = 2 + (p - 1)k + l, p \le p_{max}, 1 \le l \le k. \end{cases}$$
(9)

We can easily check that in all cases, $v_i = F_{A,v}(z_i)$ and that $\tilde{g}_i = \partial F_{A,v}(z_i)$. We then pass (v_i, \tilde{g}_i) as a response to *alg* for the query z_i , so it can update its state. Further, having defined $i_1 = 1$, the player can keep track of exploratory queries by checking whether

$$v_i \le -\frac{\eta \gamma_1}{2}$$
 and $\frac{\|P_{Span(z_{i_r}, r' \le r)^{\perp}}(z_i)\|}{\|z_i\|} \ge \frac{\gamma_2}{4}$,

where $i_1, ..., i_r$ are the indices defined so far. We perform m such iterations unless alg stops and use the last remaining queries arbitrarily. Next, we check if the last index i_k was defined. If not, we pose $i_k = m + 1$ and let z_{m+1} be the next query of alg. The final returned vectors are $z_{i_1}/\|z_{i_1}\|, ..., z_{i_k}/\|z_{i_k}\|$. This ends the description of the player's strategy.

19

20 end

```
Algorithm 1 (Strategy of the Player for the Orthogonal Vector Game with Hints)
  Input: d, k, p_{max}, m, algorithm alg
  Part 1: Strategy to store Message knowing A
      Initialize the memory of alg to be 0
      Submit \emptyset to the oracle, and use the response as v_0
      Run alg with the optimization procedure knowing A and v_0 until first exploratory query x_{i_1}
      for p \in [p_{max}], do
  5
        Let Memory<sub>v</sub> be the current memory state of alg and i_{v,1} be the current iteration step
         Run alg with the feasibility procedure until period p ends at iteration step i_{p+1,1}. If alg stopped before, return
         the strategy fails. When needed to sample a unit vector v_{p',l'}, submit vectors x_{i_{p',1}}, \dots x_{i_{p',l'}} to the oracle,
        where i_{p',1},\ldots,i_{p',r'} are the exploratory queries defined at that stage. We use the corresponding response of
         the oracle as v_{p',l'}
        if i_{p+1,1} - i_{p,1} \le m, then
           Set Message = Memory,
  8
  9 end
  10 for Remaining queries to perform to the oracle, do Submit arbitrary query (e.g., \emptyset);
  11 if Message has not been defined yet, then return The strategy fails;
  12 Submit g_{A,v} to the oracle as defined in Equation (8)
  Part 2: Strategy to make queries
  13 Set the memory state of alg to be Message, and define i_1 = 1, r = 1
  14 for i \in [m], do
  15
          Run alg with current memory to obtain a query z_i
          Submit z_i to the oracle from Game 1 to get response (g_i, s_i)
  17
          Compute (v_i, \tilde{g}_i) using z_i, g_i, and s_i as defined in Equation (9), and pass (v_i, \tilde{g}_i) as response to alg
  18
          if v_i \le -\eta \gamma_1/2 and ||P_{Span(z_{i,r}, r' \le r)^{\perp}}(z_i)||/||z_i|| \ge \gamma_2/4, then
```

Part 3: Strategy to return vectors

- 21 **if** *index* i_k *has not been defined yet*, **then**
- With the current memory of alg, find a new query z_{m+1} , and set $i_k = m + 1$
- 23 **return** $\{z_{i_1}/||z_{i_1}||, \ldots, z_{i_k}/||z_{i_k}||\}$ to the oracle

Set $i_{r+1} = i$, and increment $r \leftarrow r + 1$

We now show that the player wins with good probability. First, because alg makes at most $mp_{max} \leq d^2$ queries, by Proposition 3, on an event \mathcal{E} of probability at least $q - C\sqrt{\log d}/d$, alg succeeds and ends the period p_{max} . On \mathcal{E} , by construction, the first phase of the strategy does not fail. Now, we show that in the second phase (lines 9–12 of Game 1), the queried vectors coincide exactly with the queried vectors from the corresponding period p in the first phase (lines 3–6 of Game 1). To do so, we only need to check that the responses provided to alg coincide with the response given by the optimization procedure. First, recall that on \mathcal{E} , all periods are completed; hence, $F_{A,v,P,L} = F_{A,v}$. Next, by Proposition 1, the responses of the procedure are consistent with optimizing $F_{A,v,P,L}$ and subgradients $\partial F_{A,v,P,L}$ on an event \mathcal{F} of probability at least $1 - C'\sqrt{\log d}/d^2$. Therefore, on $\mathcal{E} \cap \mathcal{F}$, it suffices to check that the responses provided to alg are consistent with $F_{A,v}$, which we already noted; at every step i, $(v_i, \tilde{g}_i) = (F_{A,v}(z_i), \partial F_{A,v}(z_i))$. This proves that the responses and queries coincide exactly with those given by the optimization procedure on $\mathcal{E} \cap \mathcal{F}$.

Next, by construction, the chosen phase p had at most m iterations. Thus, on $\mathcal{E} \cap \mathcal{F}$, among z_1, \ldots, z_{m+1} , we have the vectors $x_{i_{p,1}}, \ldots, x_{i_{p,k}}$. Further, if i_k was not defined during part 2 of the strategy, this means that $i_k = m+1$ as defined in the player's strategy (lines 21 and 22 of Algorithm 1). As a result, for all $u \leq k$, we have $z_{i_u} = x_{i_{p,u}}$. We now show that the returned vectors $x_{i_{p,1}}/\|x_{i_{p,1}}\|, \ldots, x_{i_{p,k}}/\|x_{i_{p,k}}\|$ are successful for Game 1. First, because $i_{p,1}, \ldots, i_{p,k}$ are exploratory queries, we have directly for $u \leq k$,

$$\frac{||P_{Span(x_{i_{p,v}},v$$

Next, if l is the index of the last constructed vector $v_{p,l}$ before $i_{p,u}$ in the optimization procedure, one has $F_{A,v,p,l}(x_{i_{p,u}}) \le -\eta \gamma_1/2$. Therefore, $||Ax_{i_{p,u}}||_{\infty} \le F_{A,v,p,l}(x_{i_{p,u}}) + \eta \le \eta$. Further, $\eta v_0^{\top} x_{i_{p,u}} \le F_{A,v,p,l}(x_{i_{p,u}}) \le -\eta \gamma_1/2$.

This proves that $||x_{i_{p,u}}|| \ge \gamma_1/2$. Putting the previous two inequalities together yields

$$\frac{||Ax_{i_{p,u}}||_{\infty}}{||x_{i_{p,u}}||} \leq \frac{2\eta}{\gamma_1}.$$

As a result, this shows that the returned vectors are successful for Game 1 for the desired parameters $\alpha = 2\eta/\gamma_1$ and $\beta = \gamma_2/4$. Thus, the player wins on $\mathcal{E} \cap \mathcal{F}$, which has probability at least $q - (C + C')\sqrt{\log d}/d^2$ by the union bound. This ends the proof of the proposition. \square

3.5. Query Lower Bound for the Orthogonal Vector Game with Hints

Before proving a lower bound on the necessary number of queries for Game 1, we need to introduce two results. The first one is a known concentration result for vectors in the hypercube. It shows that for a uniform vector in the hypercube, being approximately orthogonal to k orthonormal vectors has exponentially small probability in k.

Lemma 3 (Marsden et al. [22]). Let $h \sim \mathcal{U}(\{\pm 1\}^d)$. Then, for any $t \in (0, 1/2]$ and any matrix $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_k] \in \mathbb{R}^{d \times k}$ with orthonormal columns,

$$\mathbb{P}(\|\mathbf{Z}^{\top}\boldsymbol{h}\|_{\infty} \leq t) \leq 2^{-c_H k}$$

We will also need an anticoncentration bound for random vectors, which intuitively provides a lower bound for the previous concentration result. The following lemma shows that for a uniformly random unit vector, being orthogonal to k orthonormal vectors is still achievable with exponentially small probability in k.

Lemma 4. Let k < d and $x_1, ..., x_k$ be k orthonormal vectors. Then,

$$\mathbb{P}_{\boldsymbol{y} \sim \mathcal{U}(S^{d-1})} \left(|\boldsymbol{x}_i^{\top} \boldsymbol{y}| \leq \frac{1}{d^3}, \ \forall i \leq k \right) \geq \frac{1}{e^{d-4} d^{3k}}.$$

Proof. Let $y \sim \mathcal{U}(S^{d-1})$ be a uniformly random unit vector. Then, for i < k and any y_1, \dots, y_{i-1} such that $|y_1|, \dots, |y_{i-1}| \le 1/d^3$, we have

$$\mathbb{P}\left(|y_i| \le \frac{1}{d^3} | y_1, \dots, y_{i-1}\right) = \mathbb{P}_{u \sim \mathcal{U}(S^{d-i})} \left(|u_1| \le \frac{1}{d^3 \sqrt{1 - (y_1^2 + \dots + y_{i-1}^2)}}\right) \\
\ge \frac{\int_0^{1/d^3} (1 - y^2)^{(d-i-1)/2} dy}{\int_0^1 (1 - y^2)^{(d-i-1)/2} dy} \\
\ge \frac{(1 - d^{-6})^{d/2}}{d^3} \ge \frac{e^{-d^{-5}}}{d^3},$$

where in the last equation, we used $d \ge 2$. Therefore, we can show by induction that $\mathbb{P}(|y_i| \le 1/d^3, \forall i \le k) \ge e^{-kd^{-5}}/d^{3k}$. Thus, by isometry, this shows that

$$\mathbb{P}\left(|\mathbf{x}_i^{\top}\mathbf{y}| \leq \frac{1}{d^3}, \ \forall i \leq k\right) \geq \frac{1}{e^{d-4}d^{3k}}.$$

This ends the proof of the lemma. \Box

We are now ready to prove a query lower bound for Game 1. Precisely, we show that for appropriate choices of parameters, one needs $m = \tilde{\Omega}(d)$ queries. The proof is closely inspired from the arguments given in Marsden et al. [22]. The main added difficulty arises from bounding the information leakage of the provided hints. As such, our goal is to show that these do not provide more information than the message itself.

Proposition 5. Let $k \ge 20(M + 3d \log(2d) + 1)/(c_H n)$, and let $0 < \alpha, \beta \le 1$ such that $\alpha(\sqrt{d}/\beta)^{5/4} \le 1/2$. If the player wins the orthogonal vector game with hints (Game 1) with probability at least 1/2, then $m \ge c_H d/(8(30 \log d + c_H))$.

Proof. We first define some notations. Let $Y = [y_1, \ldots, y_k]$ be the matrix storing the final outputs from the algorithm. Next, for the responses of the oracle $(g_1, s_1), \ldots, (g_m, s_m)$, we first store all of the scalar responses in a vector $c = [s_1, \ldots, s_m]$. We now focus on the responses g_1, \ldots, g_m . Next, let \tilde{G} denote the matrix containing these responses of the oracle, which are lines of A. Let G be the matrix containing unique columns from \tilde{G} augmented

with rows of A so that it has exactly m columns, which are all different rows of A. Last, let A' be the matrix A once the rows from G are removed. Next, let \tilde{V} be a matrix containing the responses of the oracle that are vectors v_l , ordered by increasing index l. As before, let V be the matrix \tilde{V} where we only conserve unique columns and append it with additional vectors v_l so that V has exactly m columns. We denote by w_1, \ldots, w_m these vectors and recall that they are vectors v_l ordered by increasing order of index l. Last, we define a vector j of indices such that j(i) contains the information of which column of the matrices G or V corresponds g_i . Precisely, if g_i is a line a from A, we set j(i) = j, where j is the index of the column from G corresponding to G. Otherwise, if G is the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G corresponding to G in the index of the column from G in the index of the column from G in the index of the column from G in the column from G in the column from G in the

Next, we argue that Y is a deterministic function of Message, the matrices G, V, and the vector of indices j and c. First, c provides the scalar responses directly. For the d-dimensional component of the responses, first note that from G, V, and j, one can easily recover the vectors g_1, \ldots, g_m . Next, using the algorithm for the second section of the orthogonal vector game with hints set with initial memory Message and the vectors g_1, \ldots, g_m as responses of the oracle, one can inductively compute the queries x_1, \ldots, x_m . Last, Y is a deterministic function of $x_i, g_i, i \in [m]$ and Message. This ends the claim that there is a function ϕ such that $Y = \phi$ (Message, G, V, j, c). Now, by the data processing inequality,

$$I(A'; Y | G, V, j, c) \le I(A'; Message | G, V, j, c) \le H(Message | G, V, j, c) \le M.$$
 (10)

In the last inequality, we used the fact that Message uses at most M bits. Now, we have that

$$I(A'; Y | G, V, j, c) = H(A' | G, V, j, c) - H(A' | Y, G, V, j, c).$$
(11)

In the next steps, we bound the two terms. We start with the second term on the right-hand side of Equation (11) using similar arguments to the proof given in Marsden et al. [22]. Let $\mathcal E$ be the event when the player succeeds at Game 1. Now, consider the case when Y is a winning matrix. Then, we have $||Ay_i||_{\infty} \leq \alpha$ for all $i \leq k$. As a result, any line a of A' satisfies $||Y^Ta||_{\infty} \leq \alpha$. Further, we have that $||P_{Span(y_j,j< i)^{\perp}}(y_i)|| \leq \beta$ for all $i \leq k$. By Lemma B.1, there exist $\lceil k/5 \rceil$ orthonormal vectors $\mathbf{Z} = [z_1, \dots, z_{\lceil k/5 \rceil}]$ such that for any $\mathbf{x} \in \mathbb{R}^d$, one has $||\mathbf{Z}^T\mathbf{x}||_{\infty} \leq (\sqrt{d}/\beta)^{5/4} ||Y^T\mathbf{x}||_{\infty}$. In particular, all lines a of A' satisfy

$$\|\mathbf{Z}^{\mathsf{T}}\mathbf{a}\|_{\infty} \leq \left(\frac{\sqrt{d}}{\beta}\right)^{5/4} \alpha \leq \frac{1}{2},$$

where we used the hypothesis in the parameters α and β . Now, by Lemma 3, one has

$$\left| \left\{ a \in \{\pm 1\}^d : \| \mathbf{Z}^{\top} a \|_{\infty} \le \frac{1}{2} \right\} \right| \le 2^d \mathbb{P}_{h \sim \mathcal{U}(\{\pm 1\}^d)} \left(\| \mathbf{Z}^{\top} h \|_{\infty} \le \frac{1}{2} \right) \le 2^{d - c_H \lceil k / 5 \rceil}.$$

Therefore, we proved that if Y' is a winning vector, $H(A'|Y=Y') \le (n-m)(d-c_Hk/5)$. Otherwise, if Y' loses, we can directly use $H(A'|Y=Y') \le (n-m)d$. Combining these equations gives

$$H(A'|Y,G,V,j,c) \leq H(A'|Y)$$

$$\leq \mathbb{P}(\mathcal{E}^c)(n-m)d + \mathbb{P}(\mathcal{E})(n-m)(d-c_Hk/5)$$

$$\leq (n-m)(d-\mathbb{P}(\mathcal{E})c_Hk/5).$$

Next, we turn to the first term of the right-hand side of Equation (11):

$$\begin{split} H(A'|G,V,j,c) &= H(A|G,V,j,c) = H(A|V) - I(A;G,j,c|V) \\ &\geq H(A|V) - H(G,j,c) \\ &\geq H(A|V) - md - m\log(2m) - m\log(d^2) \\ &= H(A) - I(A;V) - md - 3m\log(2d) \\ &= (n-m)d - 3m\log(2d) - I(A;V). \end{split}$$

In the second inequality, we use the fact that G uses md bits and j can be stored with $m \log(2m)$ bits. Now, by the chain rule,

$$I(A; V) = \sum_{i < m} I(A; w_i | w_1, \dots, w_{i-1}).$$

Now, if $w_i = v_l$, recalling that the vectors $w_{i'} = v_{l'}$ are ordered by increasing index of l', we have

$$I(A; w_{i}|w_{1},...,w_{i-1}) = H(w_{i}|w_{1},...,w_{i-1}) - H(w_{i}|A,w_{1},...,w_{i})$$

$$\leq H(w_{i}) - H(w_{i}|A,w_{1},...,w_{i},x_{l,1},...,x_{l,r_{l}})$$

$$= H(w_{i}) - H(w_{i}|x_{l,1},...,x_{l,r_{l}})$$

$$\leq \log |\mathcal{D}_{\delta}| - H(w_{i}|x_{l,1},...,x_{l,r_{l}}).$$

In the last equality, we used the fact that if $b_{l,1},\ldots,b_{l,r_l}$ are the resulting vectors from the Gram–Schmidt decomposition of $x_{l,1},\ldots,x_{l,r_l},y_l$ is generated uniformly in $S^{d-1}\cap\{y: \forall r\leq r_l, |b_{l,r}^{\top}y|\leq d^{-3}\}$ independently from the past history, and $v_l=\phi_{\delta}(y_l)$. Now, by Lemma 4, we know that

$$\mathbb{P}_{z \sim \mathcal{U}(S^{d-1})}(\forall r \leq r_l, |\boldsymbol{b}_{l,r}^{\top} \boldsymbol{z}| \leq d^{-3}) \geq \frac{1}{e^{d^{-4}} d^{3k}}.$$

As a result, for any $b_i(\delta) \in \mathcal{D}_{\delta}$, one has

$$\mathbb{P}(\boldsymbol{w}_{i} = \boldsymbol{b}_{j}(\delta) | \boldsymbol{x}_{l,1}, \dots, \boldsymbol{x}_{l,r_{l}}) \leq \frac{\mathbb{P}_{z \sim \mathcal{U}(S^{d-1})}(z \in V_{j}(\delta))}{\mathbb{P}_{z \sim \mathcal{U}(S^{d-1})}(\forall r \leq r_{l}, |\boldsymbol{b}_{l,r}^{T}\boldsymbol{z}| \leq d^{-3})} \leq \frac{e^{d^{-4}}d^{3k}}{|\mathcal{D}_{\delta}|},$$

where we used the fact that each cell has the same area. In particular, this shows that

$$H(w_i|x_{l,1},\ldots,x_{l,r_l}) = \mathbb{E}_{b \sim w_i|x_{l,1},\ldots,x_{l,r_l}}[-\log p_{w_i|x_{l,1},\ldots,x_{l,r_l}}(b)] \ge \log\left(\frac{|\mathcal{D}_{\delta}|}{e^{d^{-4}}d^{3k}}\right).$$

Hence,

$$I(A; w_i | w_1, \dots, w_{i-1}) \le 3k \log d + d^{-4} \log e.$$

Putting everything together gives

$$I(A'; Y | G, V, j) \ge (n - m)d - 3m\log(2d) - 3km\log d - 2md^{-4} - (n - m)(d - \mathbb{P}(\mathcal{E})c_H k/5)$$

$$\ge \frac{c_H}{10}k(n - m) - 3km\log d - 1 - 3d\log(2d),$$

where in the last equation, we used $d \ge 2$. Together with Equation (10), this implies

$$m \ge \frac{c_H kn/10 - M - 1 - 3d \log(2d)}{k(3\log d + c_H/10)}.$$

As a result, because $k \ge 20 \frac{M+3d \log(2d)+1}{c_H n}$ and $n \ge d/4$, we obtain

$$m \ge \frac{c_H n}{60 \log d + 2c_H} \ge \frac{c_H}{8(30 \log d + c_H)} d.$$

This ends the proof of the proposition. \Box

We are now ready to prove the main result.

Proof of Theorem 1. We set $n = \lceil d/4 \rceil$ and $k = \lceil 20(M+3d\log(2d)+1)/(c_H n) \rceil$. By Proposition 1, with probability at least $1-C\sqrt{\log d}/d^2$, the procedure is consistent with a first-order oracle for convex optimization. Hence, because the functions $F_{A,v,P,L}$ are \sqrt{d} -Lipschitz, any M-bit algorithm guaranteed to solve convex optimization within accuracy $\epsilon = \eta/(2d) = 1/d^4$ for 1-Lipschitz functions yields an algorithm that is successful for the optimization procedure with probability at least $1-C\sqrt{\log d}/d^2$ and precision $\epsilon\sqrt{d}=\eta/(2\sqrt{d})$. Suppose that it uses at most Q queries. Then, by Proposition 4, there is a strategy for Game 1 for parameters $(d,k,\lceil Q/p_{max}\rceil+1,M,\alpha=2\eta/\gamma_1,\beta=\gamma_2/4)$, in which the player wins with probability at least $1-C'\sqrt{\log d}/d$. Now, for d large enough, this probability is at least 1/2. Further,

$$\frac{2\eta}{\gamma_1} \left(\frac{4\sqrt{d}}{\gamma_2} \right)^{5/4} \le \frac{(4/3)^{5/4}}{6} \eta d^3 \le \frac{1}{2}.$$

Hence, by Proposition 5, one has

$$\lceil Q/p_{max} \rceil + 1 \ge \frac{c_H}{8(30 \log d + c_H)} d.$$

Because $p_{max} = \Theta((d/k)^{1/3} \log^{-2/3} d)$, this implies

$$Q = \Omega\left(\frac{(d/k)^{1/3}d}{\log^{5/3}d}\right) = \Omega\left(\frac{d^{5/3}}{(M + \log d)^{1/3}\log^{5/3}d}\right).$$

In particular, if $M = d^{1+\delta}$ for $\delta \in [0,1]$, the number of queries is $Q = \tilde{\Omega}(d^{1+(1-\delta)/3})$.

4. Memory-Constrained Feasibility Problem

4.1. Defining the Feasibility Procedure

Similarly to Section 3, we pose $n = \lceil d/4 \rceil$. Also, for any matrix $A \in \{\pm 1\}^{n \times d}$, we use the same functions g_A and \tilde{g}_A . We use similar techniques as those we introduced for the optimization problem. However, because in this case, the separation oracle only returns a separating hyperplane without any value considerations of an underlying function, Procedure 1 can be drastically simplified, which leads to improved lower bounds.

Let $\eta_0 = 1/(24d^2)$, $\eta_1 = 1/(2\sqrt{d})$, $\delta = 1/d^3$, and $k \le d/3 - n$ be a parameter. Last, let $p_{max} = \lfloor (c_{d,1}d - 1)/(k - 1) \rfloor$, where $c_{d,1}$ is the same quantity as in Equation (4). The feasibility procedure is defined in Procedure 2. The oracle first randomly samples $A \sim \mathcal{U}(\{\pm 1\}^{n \times d})$ and $v_0 \sim \mathcal{U}(\mathcal{D}_\delta)$. This matrix and vector are then fixed in the rest of the procedure. Whenever the player queries a point x such that $\|Ax\|_{\infty} > \eta_0$ ($v_0^\top x > -\eta_1$, respectively), the oracle returns $\tilde{g}_A(x)$ (v_0 , respectively). All other queries are called *informative* queries. With this definition, it now remains to define the separation oracle on informative queries. The oracle proceeds by periods in which the behavior is different. In each period p, the oracle constructs vectors $v_{p,1}, \ldots, v_{p,k-1}$ inductively and keeps in memory some queries $i_{p,1}, \ldots, i_{p,k}$ that will be called *exploratory*. The first informative query t will be the first exploratory query and starts period 1.

Given a new query x_t :

(*Case (f)1*) If $||Ax||_{\infty} > \eta_0$, the oracle returns $\tilde{g}_A(x_t)$.

(Case (f)2) If $v_0^{\mathsf{T}} x_t > -\eta_1$, the oracle returns v_0 .

(Case (f)3) If x_t was queried in the past sequence, the oracle returns the same vector that was returned previously.

(*Case* (f)4) Otherwise, let p be the index of the current period, and let $v_{p,1}, \ldots, v_{p,l}$ be the vectors from the current period constructed so far together with their corresponding exploratory queries $i_{p,1}, \ldots, i_{p,l} < t$. Potentially, if p = 1, one may not have defined any such vectors at the beginning of time t. In this case, let l = 0.

(Case (f)4a) If $\max_{1 \le l' \le l} v_{p,l'}^\top x_t > -\eta_1$ (with the convention $\max_{\emptyset} = -\infty$), the oracle returns $v_{p,l'}$, where $l' = \arg\max_{l \le r} v_{p,l}^\top x_t$. Ties are broken alphabetically.

(Case (f)4b) Otherwise, if l < k-1, we first define $i_{p,l+1} = t$. Then, let $\boldsymbol{b}_{p,1}, \dots, \boldsymbol{b}_{p,l+1}$ be the result from the Gram–Schmidt decomposition of $\boldsymbol{x}_{i_{p,1}}, \dots, \boldsymbol{x}_{i_{p,l+1}}$, and let $\boldsymbol{y}_{p,l+1}$ be a sample of the distribution obtained by the uniform distribution $\boldsymbol{y}_{p,l+1} \sim \mathcal{U}(S^{d-1} \cap \{\boldsymbol{z} \in \mathbb{R}^d : |\boldsymbol{b}_{p,r}^{\top}\boldsymbol{z}| \leq d^{-3}, \ \forall r \leq l+1\})$. We then pose $\boldsymbol{v}_{p,l+1} = \phi_{\delta}(\boldsymbol{y}_{p,l+1})$. Having defined this new vector, the oracle returns $\boldsymbol{v}_{p,l+1}$. We then increment l.

(Case (f)4c) Otherwise, if r=k, we define $i_{p,k}=i_{p+1,1}=t$. If $p+1\leq p_{max}$, this starts the next period p+1. As above, let $\boldsymbol{b}_{p+1,1}$ be the result of the Gram–Schmidt decomposition of $\boldsymbol{x}_{i_{p+1,1}}$ and sample $\boldsymbol{y}_{p+1,1}$ according to a uniform $\boldsymbol{y}_{p+1,1}\sim \mathcal{U}(S^{d-1}\cap\{\boldsymbol{z}\in\mathbb{R}^d:|\boldsymbol{b}_{p+1,1}^{\mathsf{T}}\boldsymbol{z}|\leq d^{-3}\})$. We then pose $\boldsymbol{v}_{p+1,1}=\phi_{\delta}(\boldsymbol{y}_{p+1,1})$, and the oracle returns $\boldsymbol{v}_{p+1,1}$. We can then increment p and reset l=1.

The above construction ends when the period p_{max} is finished. At this point, the oracle has defined the vectors $v_{p,l}$ for all $p \le p_{max}$ and $l \le k$. We then define the successful set as

$$Q_{A,v} = \left\{ x \in B_d(0,1) : ||Ax||_{\infty} \le \eta_0, v_0^{\top} x \le -\eta_1, \max_{p \le p_{max}, l \le k-1} v_{p,l}^{\top} x \le -\eta_1 \right\}.$$

From now on, the procedure uses any separation oracle for $Q_{A,v}$ as responses to the algorithm while making sure to be consistent with previous oracle responses if a query is exactly duplicated. We now define what we mean by solving the above feasibility procedure.

```
Procedure 2 (The Feasibility Procedure for Algorithm alg)
   Input: d, k, p_{max}, algorithm alg
   1 Sample A \sim \mathcal{U}(\{\pm 1\}^{n \times d}) and v_0 \sim \mathcal{U}(\mathcal{D}_\delta)
   2 Initialize the memory of alg to 0, and let p = 1, l = 0
   3 for t \ge 1, do
           Run alg with current memory to obtain a query x_t
           if ||Ax_t|| > \eta_0, then return \tilde{g}_A(x_t) as response to alg;
   6
           else if v_0^{\mathsf{T}} x_t > -\eta_1, then return v_0 as response to alg;
   7
           else if Query x_t was made in the past, then return same vector that was returned for x_t;
   8 else
          \begin{aligned} \text{if } \max_{1 \leq l' \leq l} v_{p,l'}^\top x_t > -\eta_1, \text{then} \\ \text{return } v_{p,l'}, \text{ where } l' = \text{arg } \max_{l \leq r} v_{p,l}^\top x_t \end{aligned}
   9
   10
   11
                    else if l < k - 1, then
                        Let i_{p,l+1} = t, and compute Gram–Schmidt decomposition \boldsymbol{b}_{p,1}, \dots, \boldsymbol{b}_{p,l+1} of \boldsymbol{x}_{i_{p,1}}, \dots, \boldsymbol{x}_{i_{p,l+1}} Sample \boldsymbol{y}_{p,l+1} uniformly on \mathcal{S}^{d-1} \cap \{\boldsymbol{z} \in \mathbb{R}^d : |\boldsymbol{b}_{p,l'}^\top \boldsymbol{z}| \leq d^{-3}, \ \forall l' \leq l+1\}, and define
   12
   13
   14
   15
                        return v_{p,l+1} as response to alg and increment l \leftarrow l+1
                    else if p + 1 \le p_{max}, then
                        Set i_{p,k} = i_{p+1,1} = t, and compute the Gram–Schmidt decomposition b_{p+1,1} of x_{i_{p+1,1}}
   17
                        Sample y_{p+1,1} uniformly on S^{d-1} \cap \{z \in \mathbb{R}^d : |b_{p+1,1}^\top z| \le d^{-3}\}, and define v_{p+1,1} = \phi_{\delta}(y_{p+1,1})
   18
                        return v_{p+1,1} as response to alg, increment p \leftarrow p+1, and reset l=1
   19
                    else Set i_{p_{max},k} = t, and break the for loop;
   20
   22 for t' \ge t, do Use any separation oracle for Q_{A,v} consistent with previous responses
```

Definition 3. Let alg be an algorithm for the feasibility problem. When running alg with the responses of the feasibility procedure, we denote by v the set of constructed vectors and $x^*(alg)$ the final answer returned by alg. We say that an algorithm alg is successful for the feasibility procedure with probability $q \in [0,1]$ if taking $A \sim \mathcal{U}(\{\pm 1\}^{n\times d})$ with probability at least q over the randomness of A and of the procedure, $x^*(alg) \in Q_{A,v}$.

In the rest of this section, we first relate this feasibility procedure to the standard feasibility problem, and then, we prove query lower bounds to solve the feasibility procedure.

4.2. Reduction from the Feasibility Procedure to the Feasibility Problem

In the next proposition, we check that the above procedure indeed corresponds to a valid feasibility problem.

Proposition 6. On an event of probability at least $1 - C\sqrt{\log d}/d$, the procedure described above is a valid feasibility problem. More precisely, the following hold.

• There exists $\bar{x} \in B_d(0,1)$ such that $||A\bar{x}||_{\infty} = 0$, $v_0^{\top}\bar{x} \leq -4\eta_1$, and

$$\max_{p \leq p_{\max}, l \leq k-1} v_{p,l}^{\top} \bar{x} \leq -4\eta_1.$$

- Let $\epsilon = \min\{\eta_0/\sqrt{d}, \eta_1\}/2$. Then, $B_d(\bar{x} \epsilon(\bar{x}/||\bar{x}||), \epsilon) \subseteq B_d(0, 1) \cap B_d(\bar{x}, 2\epsilon) \subseteq Q_{A,v}$.
- Throughout the run of the feasibility problem, the separation oracle always returned a valid cut: that is, for any iteration t, if x_t denotes the query and g_t is the returned vector from the oracle, one has

$$\forall x \in Q_{A,v}, \quad \langle g_t, x_t - x \rangle > 0.$$

Further, responses are consistent; if $x_t = x_{t'}$, the responses of the procedure at times t and t' coincide.

We use a similar proof to that of Proposition 2.

Proof. For convenience, we rename $v_{p,l} = v_{(p-1)(k-1)+l}$. Also, let $l_{max} = p_{max}(k-1) \le c_{d,1}d - 1$. Next, let $C_d = \sqrt{40l_{max}\log d}$. We define the vector

$$\bar{x} = -\frac{1}{C_d} \sum_{l=0}^{l_{max}} P_{Span(a_i, i \le n)^{\perp}}(v_l).$$

Because $l_{max} \le p_{max}(k-1) \le c_{d,1}d-1$, the same arguments as in the proof of Proposition 2 show that on an event \mathcal{E} of probability at least $1 - C\sqrt{\log d}/d$, we have $\|\bar{x}\| \le 1$ and

$$\max_{0 \leq l \leq l_{max}} \boldsymbol{v}_l^\top \bar{\boldsymbol{x}} \leq -\frac{1}{40\sqrt{(l_{max}+1)\mathrm{log}\,d}} \leq -\frac{2}{\sqrt{d}} = -4\eta_1,$$

where in the second inequality, we used $l_{max} \leq c_{d,1}d-1$. Now, by construction, one has $||A\bar{x}||_{\infty} = 0$. This ends the proof of the first claim of the proposition. We now turn to the second claim, which is immediate from the fact that $x \mapsto ||Ax||_{\infty}$ is \sqrt{d} -Lipschitz and both $x \mapsto v_0^{\top}x$ and $x \mapsto \max_{p \leq p_{max}, l \leq k} v_{p,l}^{\top}x$ are 1-Lipschitz. Therefore, $B_d(\bar{x} - \epsilon \bar{x}/||\bar{x}||, \epsilon) \subseteq B_d(0, 1) \cap B_d(\bar{x}, 2\epsilon) \subset Q_{A,v}$. It now remains to check that the third claim is satisfied. It suffices to check that this is the case during the construction phase of the feasibility procedure: by construction of $Q_{A,v} \subset \{x : ||Ax||_{\infty} \leq \eta_0\}$.

Hence, it suffices to check that for informative queries x_t , the returned vectors g_t are valid separating hyperplanes. By construction, these can only be either v_0 or $v_{p,l}$ for $p \le p_{max}$, $l \le k-1$. We denote by w this vector. Let t' be the first time x_t was queried. There are two cases. Either w was not constructed at time t', in which case, by construction this means that we are in case ((f)2) or case ((f)4a). Both cases imply $w^Tx_t > -\eta_1$. Hence, w, which is returned by the procedure, is a valid separating hyperplane. Now, suppose that $w = v_{p,l}$ was constructed at time t'—case ((f)4b) or case ((f)4c). By construction, one has $|b_{p,r}^Ty_{p,l}| \le d^{-3}$ for all $r \le l$. Decomposing $x_t = x_{l_{p,l}} = \alpha b_{p,1} + \cdots + \alpha_l b_{p,l}$, we obtain

$$|x_t^{\mathsf{T}} y_{p,l}| \le \frac{\|\alpha\|_1}{d^3} \le \frac{1}{d^2 \sqrt{d}}.$$

As a result, $y_{p,l}^{\top} x_t \ge -1/(d^2 \sqrt{d})$. Now, because $v_{p,l} = \phi_{\delta}(y_{p,l})$, we have $||v_{p,l} - y_{p,l}|| \le \delta$. Hence, for any $d \ge 2$,

$$\boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_t \geq -1/(d^2 \sqrt{d}) - \delta > -\eta_1.$$

Hence, w was a valid separating hyperplane. The last claim that the responses of the procedure are consistent over time is a direct consequence from its construction. This ends the proof of the proposition. \Box

As a simple consequence of this result, solving the feasibility problem is harder than solving the feasibility procedure with high probability.

Proposition 7. Let alg be an algorithm that solves the feasibility problem with accuracy $\epsilon = 1/(48d^2\sqrt{d})$. Then, it solves the feasibility procedure with probability at least $1 - C\sqrt{\log d}/d$.

Proof. Let \mathcal{E} be the event of probability at least $1 - C\sqrt{\log d}/d$ defined in Proposition 6. We show that on \mathcal{E} , alg solves the feasibility procedure. On \mathcal{E} , the feasibility procedure emulates a valid feasibility oracle. Further, on \mathcal{E} , the successful set contains a closed ball of radius ε . As a result, on \mathcal{E} , alg finds a solution to the feasibility problem emulated by the procedure. \square

Next, we show that it is necessary to finish the p_{max} periods to solve the feasibility procedure.

Proposition 8. Fix an algorithm alg. Then, if \mathcal{E} denotes the event when alg succeeds and \mathcal{B} denotes the event when the procedure ends period p_{max} with alg, then $\mathcal{E} \subseteq \mathcal{B}$.

Proof. Consider the case when the period p_{max} was not ended. Let x^* denote the last query performed by alg. We consider the scenario in which x^* fell. Let t be the first time when alg submitted query x^* . For any of case ((f)1), case ((f)2), or case ((f)4a), by construction of $Q_{A,v}$, we already have $x_t \notin Q_{A,v}$. It remains to check case ((f)4b) and case ((f)4c), for which the procedure constructs a new vector $v_{p,l}$, where p is the index of the period of t and $i_{p,1}, \ldots, i_{p,l} = t$ are the previous exploratory queries in period p. We decompose $x_t = x_{i_{p,l}} = \alpha_1 b_{p,1} + \alpha_l b_{p,l}$. Now, by construction,

$$|\mathbf{x}_{t}^{\mathsf{T}}\mathbf{y}_{p,l}| = |\mathbf{x}_{i_{p,l}}^{\mathsf{T}}\mathbf{y}_{p,l}| \le \frac{\|\mathbf{\alpha}\|_{1}}{d^{3}} \le \frac{1}{d^{2}\sqrt{d}}.$$

As a result, $\mathbf{x}_t^{\mathsf{T}} \mathbf{v}_{p,l} \ge -|\mathbf{x}_t^{\mathsf{T}} \mathbf{y}_{p,l}| - \delta \ge -d^{-2.5} - d^{-3} > -\eta_1$ for any $d \ge 2$. Thus, $\mathbf{x}_t = \mathbf{x}^{\star} \notin Q_{A,v}$. This shows that in order to succeed at the feasibility procedure, an algorithm needs to end all p_{max} periods. \square

4.3. Reduction from the Orthogonal Vector Game with Hints

The remaining piece of our argument is to show that solving the feasibility procedure is harder than solving the orthogonal vector game with hints: Game 1.

Proposition 9. Let $A \sim \mathcal{U}(\{\pm 1\}^{n \times d})$. If there exists an M-bit algorithm that solves the feasibility problem described above using mp_{max} queries with probability at least q over the randomness of the algorithm, choice of A, and the randomness of the separation oracle, then there is an algorithm for Game 1 for parameters $(d,k,m,M,\alpha=\eta_0/\eta_1,\beta=\eta_1/2)$, for which the player wins with probability at least q over the randomness of the player's strategy and A.

Proof. Let alg be an M-bit algorithm solving the feasibility problem with mp_{max} queries with probability at least q. In Algorithm 2, we describe the strategy of the player in Game 1.

Algorithm 2 (Strategy of the Player for the Orthogonal Vector Game with Hints)

Input: d, k, p_{max} , m, algorithm alg

Part 1: Strategy to store Message knowing *A*

- 1 Initialize the memory of *alg* to be **0**
- 2 Submit \emptyset to the oracle, and use the response as v_0
- 3 Run alg with the optimization procedure knowing A and v_0 until the first exploratory query $x_{i_{1,1}}$
- 4 for $p \in [p_{max}]$, do
- Let Memory, be the current memory state of alg and $i_{p,1}$ be the current iteration step
- Run alg with the feasibility procedure until period p ends at iteration step $i_{p+1,1}$. If alg stopped before, **return** the strategy fails. When needed to sample a unit vector $v_{p',l'}$, submit vectors $x_{i_{p',1}}$, ... $x_{i_{p',l'}}$ to the oracle. We use the corresponding response of the oracle as $v_{p',l'}$
- 7 if $i_{p+1,1} i_{p,1} \le m$, then
- 8 Set Message = Memory_p
- 9 end
- 10 **for** *Remaining queries to perform to the oracle,* **do** Submit arbitrary query (e.g., \emptyset).
- 11 **if** Message has not been defined yet, **then return** The strategy fails;
- 12 Submit $\tilde{g}_{A,v}$ to the oracle as defined in Equation (12)

Part 2: Strategy to make queries

- 13 Set the memory state of alg to be Message
- 14 for $i \in [m]$, do
- Run *alg* with current memory to obtain a query z_i
- Submit z_i to the oracle from Game 1 to get response (g_i, s_i)
- 17 Compute \tilde{g}_i using z_i , g_i , and s_i as defined in Equation (13), and pass \tilde{g}_i as response to alg
- 18 **end**

Part 3: Strategy to return vectors

- 19 **for** $l \in [k]$, **do** Set i_l to be the index i of the first query z_i for which $s_i = l$ if it exists;
- 20 **if** $index i_k$ has not been defined yet, **then**
- 21 With the current memory of alg, find a new query z_{m+1} , and set $i_k = m + 1$
- 22 **return** $\{z_{i_1}/||z_{i_1}||, \ldots, z_{i_k}/||z_{i_k}||\}$ to the oracle

In the first part of the strategy, the player observes A. Then, they proceed to simulate the feasibility problem with alg using parameters A. When needed to sample a vector $v_{p,l}$ (v_0 , respectively), the player submits the corresponding queries $x_{i_{p,1}}, \ldots, x_{i_{p,l}}$ (\emptyset , respectively) useful to define $v_{p,l}$. The player then takes the response given by the oracle as that vector $v_{p,l}$ (v_0 , respectively), which simulates exactly a run of the feasibility procedure. Further, because $1 + p_{max}(k-1) \le d$, the player does not run out of queries. Importantly, during the run, the player keeps track of the length $i_{p,k} - i_{p,1}$ of period p. The first time we encounter a period p with length at most m, we set Message = Memory $_p$, the memory state of alg at the beginning of period p. If there is no such period, the strategy fails. Also, if alg stopped before ending period p_{max} , the strategy fails. Next, the algorithm submits the following function $\tilde{g}_{A,v}$ to the oracle. Because the responses of the feasibility procedure are consistent over time, we adopt the following notation. For a previously queried vector x of alg, we denote g(x) the vector, which was returned to alg during the first part (lines 3–9 of Algorithm 2):

$$\tilde{g}_{A,v}: x \longmapsto
\begin{cases}
(0,1) & \text{if } x \text{ was never queried in the first part,} \\
(a_{i},1) & \text{otherwise and if } g(x) \in \{\pm a_{i}\}, i \leq n, \\
(v_{0},2) & \text{otherwise and if } g(x) = v_{0}, \\
(v_{p',l'},2+l'\mathbb{1}_{p'=p}+k\mathbb{1}_{p'=p+1,l'=1}) & \text{otherwise and if } g(x) = v_{p',l'}, p' \leq p_{\max}, 1 \leq k-1.
\end{cases} (12)$$

Intuitively, the first component of \tilde{g} gives the returned vector in the first period at the exception that we always return a_i instead of $\{\pm a_i\}$. The second term has values in $[2+k \le d^2]$. Hence, the submitted function is valid.

Next, in the second part of the algorithm, the player proceeds to simulate a run of the feasibility procedure with alg on period p. To do so, we first set the memory state of alg to Message. Each new query z_i is submitted to the oracle of Game 1 to get a response (g_i, s_i) . Then, we compute \tilde{g}_i as follows:

$$\tilde{\mathbf{g}}_i = \begin{cases} \mathbf{g}_i & \text{if } s_i \ge 2, \\ sign(\mathbf{g}_i^{\top} \mathbf{z}_i) \mathbf{g}_i & \text{if } s_i = 1. \end{cases}$$
 (13)

One can easily check that \tilde{g}_i corresponds exactly to the response that was passed to alg in the first part of the strategy. The player then passes \tilde{g}_i to alg so that it can update its state. We repeat this process for m steps. Further, the player can also keep track of the exploratory queries; the index i_l of the first response satisfying $s_i = 2 + l$ for $l \leq k-1$ ($s_i = 2 + k$, respectively) is the exploratory query, which led to the construction of $v_{p,l}$ ($v_{p+1,1}$, respectively) in the first part. Last, we check if the last index i_k was defined. If not, we pose $i_k = m+1$ and let z_{m+1} be the next query of alg with the current memory. The player then returns the vectors $z_{i_1}/\|z_{i_1}\|,\ldots,z_{i_k}/\|z_{i_k}\|$. This ends the description of the player's strategy.

By Proposition 8, on an event \mathcal{E} of probability at least q, the algorithm alg succeeds and ends period p_{max} . As a result, similarly as in the proof of Proposition 4, because alg makes at most mp_{max} queries and there are p_{max} periods, there must be a period of length at most m. Hence, the strategy never fails at this phase of the player's strategy on the event \mathcal{E} . Further, we already checked that in the second phase, the vectors \tilde{g}_i passed to alg coincide exactly with the responses passed to alg in the first part. Thus, this shows that during the second part, the player simulates exactly the run of the feasibility problem on period p. More precisely, the queries coincide with the queries in the feasibility problem at times $i_{p,1},\ldots, \min\{i_{p,k},i_{p,1}+m-1\}$. Now, because the first part succeeded on \mathcal{E} , we have $i_{p,k} \leq i_{p,0}+m$. Therefore, if i_k has not yet been defined, this means that we had $i_{p,k}=i_{p,1}+m$. Hence, the next query with the current memory z_{m+1} is exactly the query $x_{i_{p,k}}$ for the feasibility problem. This shows that the vectors z_{i_1},\ldots,z_{i_k} coincide exactly with the vectors $x_{i_{p,1}},\ldots,x_{i_{p,k}}$ when running alg on the feasibility problem in the first part.

We now show that the returned vectors are successful for Game 1. By construction, $x_{i_{p,1}}, \ldots, x_{i_{p,k}}$ are all informative. In particular, $||Ax_{i_{p,l}}||_{\infty} \leq \eta_0$ for all $1 \leq l \leq k$. Further, these queries did not fall in case ((f)2); hence, $v_0^{\mathsf{T}}x_{i_{p,l}} < -\eta_1$, which implies $||x_{i_{p,l}}|| > \eta_1$ for all $l \leq k$. As a result,

$$\frac{\|Ax_{i_{p,l}}\|_{\infty}}{\|x_{i_{p,l}}\|} \leq \frac{\eta_0}{\eta_1}.$$

Next, fix $l \le k-1$. By construction of $y_{n,l}$,

$$||P_{Span(\mathbf{x}_{i_{p,l'}},l'\leq l)}(\mathbf{y}_{p,l})||^2 = \sum_{l'\leq l} |\mathbf{b}_{p,l'}^{\top} \mathbf{y}_{p,l}|^2 \leq \frac{k}{d^6} \leq \frac{1}{d^5}.$$

Hence,

$$||v_{p,l} - P_{Span(x_{i_{p,l'}}, l' \le l)^{\perp}}(y_{p,l})|| \le ||P_{Span(x_{i_{p,l'}}, l' \le l)}(y_{p,l})|| + \delta \le \frac{1}{d^5} + \delta.$$

As a result, because $x_{p,l+1}^{\top}v_{p,l}<-\eta_1$, we have

$$||P_{Span(\mathbf{x}_{i_{v,l'}}, l' \leq l)^{\perp}}(\mathbf{x}_{p, l+1})|| \geq |\mathbf{x}_{p, l+1}^{\top} P_{Span(\mathbf{x}_{i_{v,l'}}, l' \leq l)^{\perp}}(\mathbf{y}_{p, l})| > \eta_1 - \frac{1}{d^5} - \delta \geq \frac{\eta_1}{2}.$$

This shows that the returned vectors $\mathbf{x}_{i_p,1}/\|\mathbf{x}_{i_p,1}\|,\ldots,\mathbf{x}_{i_p,k}/\|\mathbf{x}_{i_p,k}\|$ are successful for Game 1 with parameters $\alpha=\eta_0/\eta_1$ and $\beta=\eta_1/2$. This ends the proof that strategy succeeds on \mathcal{E} for these parameters, which ends the proof of the proposition. \square

We are now ready to prove the main result.

Proof of Theorem 2. Suppose that there is an algorithm *alg* for solving the feasibility problem to optimality $\epsilon = 1/(48d^2\sqrt{d})$ with memory M and at most Q queries. Let $k = \lceil 20(M+3d\log(2d)+1)/(c_Hn) \rceil$. By Proposition 7, it solves the feasibility procedure with parameter k with probability at least $1 - C\sqrt{\log d}/d$. By Proposition 9, there is an algorithm for Game 1 that wins with probability 1/3 with $m = \lceil Q/p_{max} \rceil$ and parameters $\alpha = \eta_0/\eta_1$ and $\beta = \eta_1/2$. Now, we check that

$$\alpha \left(\frac{\sqrt{d}}{\beta}\right)^{5/4} \le 12d^2\eta_0 = \frac{1}{2}.$$

Hence, by Proposition 5, we have

$$m \ge \frac{c_H}{8(30 \log d + c_H)} d.$$

This shows that

$$Q \geq \Omega\left(p_{max}\frac{d}{\log d}\right) = \Omega\left(\frac{d^2}{k\log^3 d}\right) = \Omega\left(\frac{d^3}{(M + \log d)\log^3 d}\right).$$

This implies that for a memory $M = d^{2-\delta}$ with $0 \le \delta \le 1$, the number of queries is $Q = \tilde{\Omega}(d^{1+\delta})$. \square

5. Conclusion and Future Directions

In this work, we established lower bounds for the query complexity of memory-constrained algorithms for convex optimization and its related feasibility problem. Our findings highlight that quadratic memory is necessary for achieving the optimal oracle complexity in first-order convex optimization. By establishing these lower-bound trade-offs, our research contributes to a deeper understanding of the computational aspects of convex optimization.

It is worth noting that our lower bounds only apply to deterministic algorithms. Although many standard optimization methods are deterministic, generalizing our results to randomized algorithms is also desirable. We note that subsequent to our work, Chen and Peng [10] gave slightly weaker lower bounds, which hold for randomized algorithms and up to near-quadratic memory as well.

Last, providing memory-constrained algorithms for convex optimization beyond the standard cutting-plane methods and gradient descent approaches is an important question. As depicted in Figure 1 (see Marsden et al. [22] and Woodworth and Srebro [42]), to the best of our knowledge, no algorithms from the literature provided such oracle complexity/memory trade-offs in any regime $\epsilon \ll 1/\sqrt{d}$. The authors are investigating this question, and in a recent follow-up (Blanchard et al. [6]), we proposed a family of memory-constrained algorithms parametrized by $p \in [d]$, which provides an oracle complexity/memory trade-off for subpolynomial regimes: $\ln \frac{1}{\epsilon} \gg \ln d$. Importantly, in the exponential regime $\epsilon \leq d^{-\Omega(d)}$, our algorithm with p = d improves the oracle complexity of gradient descent while preserving the same memory usage.

Acknowledgments

A previous version of this work appeared at the 36th Annual Conference on Learning Theory.

Appendix A. Concentration Bounds

The following result gives concentration bounds for the norm of the projection of a random unit vector onto linear subspaces.

Proposition A.1. Let P be a projection in \mathbb{R}^d of rank r, and let $x \in \mathbb{R}^d$ be a random vector sampled uniformly on the unit sphere $x \sim \mathcal{U}(S^{d-1})$. Then, for every t > 0,

$$\max\left\{\mathbb{P}\left(\|P(x)\|^2 - \frac{r}{d} \ge t\right), \mathbb{P}\left(\|P(x)\|^2 - \frac{r}{d} \le -t\right)\right\} \le e^{-dt^2}.$$

Further, if r = 1 and $d \ge 2$,

$$\mathbb{P}\left(\|P(x)\| \ge \sqrt{\frac{t}{d-1}}\right) \le 2\sqrt{t}e^{-t/2}.$$

Proof. First, by isometry, we can assume that P is the projection onto the coordinate vectors $e_1, \dots e_r$. Then, let $y \sim \mathcal{N}(0,1)$ be a normal vector. Note that $x = y/\|y\| \sim \mathcal{U}(S^{d-1})$. Further,

$$||x||^2 \ge \frac{r}{d} + t \Longleftrightarrow \left(1 - \frac{r}{d} - t\right) \sum_{i=1}^r y_i^2 \ge \left(\frac{r}{d} + t\right) \sum_{i=r+1}^d y_i^2.$$

Now, note that $Z_1 = \sum_{i=1}^r y_i^2$ and $Z_2 = \sum_{i=r+1}^d y_i^2$ are two independent random chi-squared variables of parameters r and d - r, respectively. Recall that the moment-generating function of $Z \sim \chi^2(k)$ is $\mathbb{E}[e^{sZ}] = (1-2s)^{-k/2}$ for s < 1/2. Therefore, for any

$$-\frac{1}{2(r/d+t)} < s < \frac{1}{2(1-r/d-t)},\tag{A.1}$$

one has

$$\mathbb{P}\left(\|P(x)\|^{2} - \frac{r}{d} \ge t\right) \le \mathbb{E}\left[\exp\left(s\left(1 - \frac{r}{d} - t\right)Z_{1} - s\left(\frac{r}{d} + t\right)Z_{2}\right)\right]$$

$$= \frac{\left[1 - 2s\left(1 - \frac{r}{d} - t\right)\right]^{-r/2}}{\left[1 + 2s\left(\frac{r}{d} + t\right)\right]^{-(d-r)/2}}.$$

Now, let

$$s = \frac{1}{2} \left(\frac{1 - r/d}{1 - r/d - t} - \frac{r/d}{r/d + t} \right),$$

which satisfies Equation (A.1). The previous equation readily yields

$$\mathbb{P}\left(\left\|P(x)\right\|^2 - \frac{r}{d} \ge t\right) \le \exp\left(-\frac{d}{2}d_{KL}\left(\frac{r}{d}; \frac{r}{d} + t\right)\right) \le e^{-dt^2}.$$

In the last inequality, we used Pinsker's inequality $d_{KL}(r/d;r/d+t) \ge 2\delta(\mathcal{B}(r/d),\mathcal{B}(d/r+t))^2 = 2t^2$, where $\mathcal{B}(q)$ is the Bernoulli distribution of parameter q. Replacing P with Id-P and r with d-r gives the other inequality:

$$\mathbb{P}\Big(\|P(x)\|^2 - \frac{r}{d} \le -t\Big) \le e^{-dt^2}.$$

This gives the first claim. For the second claim, supposing that r = 1 < d, from the above equation, we have

$$\mathbb{P}\left(\|P(x)\|^{2} \ge \frac{t}{d}\right) \le \exp\left(-\frac{d}{2}d_{KL}\left(\frac{1}{d}; \frac{t}{d}\right)\right) = \sqrt{t}\left(\frac{1-\frac{t}{d}}{1-\frac{1}{d}}\right)^{(d-1)/2} \le \sqrt{2t}e^{-t(d-1)/(2d)}.$$

Thus,

$$\mathbb{P}\left(\|P(x)\|^2 \ge \frac{t}{d-1}\right) \le \sqrt{\frac{2(d-1)}{d}}\sqrt{t}e^{-t/2},$$

which ends the proof of the proposition. \Box

Next, we need the following lemma, which gives a concentration inequality for discretized samples in \mathcal{D}_d and is approximately perpendicular to $k \le d/3 - 1$ vectors.

Lemma A.1. Let $0 \le k \le d/3 - 1$ and $x_1, \ldots, x_k \in B_d(0,1)$ be k orthonormal vectors in the unit ball, and $\mathbf{x} \in B_d(0,1)$. Denote by μ the distribution on the unit sphere corresponding to the uniform distribution $\mathbf{y} \sim \mathcal{U}(S^{d-1} \cap \{\mathbf{w} \in \mathbb{R}^d : |\mathbf{x}_i^\top \mathbf{w}| \le d^{-3}, \ \forall i \le k\})$. Let $\mathbf{y} \sim \mu$. Then, for $t \ge 2$,

$$\mathbb{P}\left(|\mathbf{x}^{\top}\mathbf{y}| \ge \sqrt{\frac{t}{d}} + \frac{1}{d^2}\right) \le 2\sqrt{t}e^{-t/3}.$$

Further, let $\delta \leq 1$ and $z = \phi_{\delta}(y)$. Then, for $t \geq 4$,

$$\mathbb{P}\left(|x^{\top}z| \ge \sqrt{\frac{t}{d}} + \frac{1}{d^2} + \delta\right) \le 2\sqrt{t}e^{-t/3}.$$

Proof. We use the same notations as above and denote by $\mathcal{E} = \{|x_i^\top y| \le d^{-3}, \ \forall i \le k\}$ the event considered and $y \sim \mu$. We decompose $y = \alpha_1 x_1 + \dots + \alpha_k x_k + y'$, where $y' \in Span(x_i, i \le k)^\perp := E$. Now, note that $y'/\|y'\|$ is a uniformly random unit vector in E. As a result, using Proposition A.1, we obtain for any $t \ge 2$,

$$\mathbb{P}\left(|x^{\top}y'| \ge \sqrt{\frac{t}{d-k-1}}\right) = \mathbb{P}\left(|P_E(x)^{\top}y'| \ge \sqrt{\frac{t}{d-k-1}}\right)$$

$$< 2\sqrt{t}e^{-t/2}.$$

Also, because by definition of μ , we have $|\alpha_i| \le d^{-3}$ for all $i \le k$, we obtain $|x^T y| \le k/d^3 + |x^T y'| \le 1/d^2 + |x^T y'|$. As a result, using the fact that $d - k - 1 \ge 2d/3$, the previous equation shows that

$$\mathbb{P}\left(|x^\top y| \geq \sqrt{\frac{3t}{2d}} + \frac{1}{d^2}\right) \leq \mathbb{P}\left(|x^\top y'| \geq \sqrt{\frac{t}{d-k-1}}\right) \leq 2\sqrt{t}e^{-t/2}.$$

Next, we use the fact that $||z-y|| = ||\phi_{\delta}(y)-y|| \le \delta$ to obtain

$$\mathbb{P}\left(|x^{\top}z| \geq \sqrt{\frac{t}{d}} + \frac{1}{d^2} + \delta\right) \leq \mathbb{P}\left(|x^{\top}y| \geq \sqrt{\frac{t}{d}} + \frac{1}{d^2}\right) \leq 2\sqrt{t}e^{-t/3}.$$

This ends the proof of the lemma. \Box

Appendix B. An Improved Result on Robustly Independent Vectors

The following lemma serves the same purpose as Marsden et al. [22, lemma 34]. Namely, from successful vectors of Game 1, it allows us to recover an orthonormal basis that is still approximately in the null space of A. The following version gives a stronger version that improves the dependence in d of our chosen parameters.

Lemma B.1. Let $\delta \in (0,1]$, and suppose that we have $r \leq d$ unit norm vectors $y_1, \ldots, y_r \in \mathbb{R}^d$. Suppose that for any $i \leq k$,

$$||P_{Span(\boldsymbol{y}_i,j< i)^{\perp}}(\boldsymbol{y}_i)|| \geq \delta.$$

Let $Y = [y_1, \dots, y_r]$ and $s \ge 2$. There exists $\lceil r/s \rceil$ orthonormal vectors $\mathbf{Z} = [z_1, \dots, z_{\lceil r/s \rceil}]$ such that for any $\mathbf{a} \in \mathbb{R}^d$,

$$\|\mathbf{Z}^{\mathsf{T}}\mathbf{a}\|_{\infty} \leq \left(\frac{\sqrt{d}}{\delta}\right)^{s/(s-1)} \|\mathbf{Y}^{\mathsf{T}}\mathbf{a}\|_{\infty}.$$

Proof. Let $B = (b_1, ..., b_r)$ be the orthonormal basis given by the Gram–Schmidt decomposition of $y_1, ..., y_r$. By definition of the Gram–Schmidt decomposition, we can write Y = BC, where C is an upper-triangular matrix. Further, its diagonal is exactly $diag(\|P_{Span(y_r, l' < l)^{\perp}}(y_l)\|, l \le r)$. Hence,

$$\det(\mathbf{Y}) = \det(\mathbf{C}) = \prod_{l < r} ||P_{Span(\mathbf{y}_{l'}, l' < l)^{\perp}}(\mathbf{y}_{l})|| \ge \delta^{r}.$$

We now introduce the singular value decomposition $Y = Udiag(\sigma_1, ..., \sigma_r)V^{\top}$, where $U \in \mathbb{R}^{d \times r}$ and $V \in \mathbb{R}^{r \times r}$ have orthonormal columns and $\sigma_1 \ge \cdots \ge \sigma_r$. Next, for any vector $z \in \mathbb{R}^d$, because the columns of Y have unit norm,

$$||Yz||_2 \le \sum_{l \le r} |z_l| ||y_l||_2 \le ||z||_1 \le \sqrt{d} ||z||_2.$$

In the last inequality, we used Cauchy–Schwartz. Therefore, all singular values of Y are upper bounded by $\sigma_1 \le \sqrt{d}$. Thus, with $r' = \lceil r/s \rceil$,

$$\delta^{r} \le \det(Y) = \prod_{l=1}^{r} \sigma_{l} \le d^{(r'-1)/2} \sigma_{r'}^{r-r'+1} \le d^{r/2s} \sigma_{r'}^{(s-1)r/s}$$

so that $\sigma_{r'} \ge \delta^{s/(s-1)}/d^{1/(2s)}$. We are ready to define the new vectors. We pose for all $i \le r'$, $z_i = u_i$ the ith column of U. These correspond to the r' largest singular values of Y and are orthonormal by construction. Then, for any $i \le r'$, we also have $z_i = u_i = Yv_i/\sigma_i$, where v_i is the ith column of V. Hence, for any $a \in \mathbb{R}^d$,

$$|z_i^{\top} a| = \frac{1}{\sigma_i} |v_i^{\top} Y^{\top} a| \le \frac{\|v_i\|_1}{\sigma_i} \|Y^{\top} a\|_{\infty} \le \frac{d^{1/2 + 1/(2s)}}{\delta^{s/(s - 1)}} \|Y^{\top} a\|_{\infty}.$$

This ends the proof of the lemma. \Box

Endnote

¹ $\tilde{\Omega}$ and \tilde{O} hide polylog(*d*) factors.

References

- [1] Anstreicher KM (2000) The volumetric barrier for semidefinite programming. Math. Oper. Res. 25(3):365–380.
- [2] Atkinson DS, Vaidya PM (1995) A cutting plane algorithm for convex programming that uses analytic centers. *Math. Programming* 69(1–3):1–43.
- [3] Balkanski E, Singer Y (2018) Parallelization does not accelerate convex optimization: Adaptivity lower bounds for non-smooth convex minimization. Preprint, submitted August 12, https://arxiv.org/abs/1808.03880.
- [4] Beame P, Oveis Gharan S, Yang X (2018) Time-space tradeoffs for learning finite functions from random evaluations, with applications to polynomials. Bubeck S, Perchet V, Rigollet P, eds. *Proc. 31st Conf. Learn. Theory*, vol. 75 (PMLR, New York), 843–856.
- [5] Bertsimas D, Vempala S (2004) Solving convex programs by random walks. J. ACM 51(4):540-556.
- [6] Blanchard M, Zhang J, Jaillet P (2023) Memory-constrained algorithms for convex optimization. Oh A, Naumann T, Globerson A, Saenko K, Hardt M, Levine S, eds. *Advances in Neural Information Processing Systems*, vol. 36 (Curran Associates Inc., Red Hook, NY), 6156–6189.
- [7] Brown G, Bun M, Smith A (2022) Strong memory lower bounds for learning natural models. Loh P-L, Raginsky M, eds. *Proc. Thirty Fifth Conf. Learn. Theory*, vol. 178 (PMLR, New York), 4989–5029.
- [8] Brown G, Bun M, Feldman V, Smith A, Talwar K (2021) When is memorization of irrelevant training data necessary for high-accuracy learning? Proc. 53rd Annual ACM SIGACT Sympos. Theory Comput. (Association for Computing Machinery, New York), 123–132.
- [9] Bubeck S, Jiang Q, Lee YT, Li Y, Sidford A (2019) Complexity of highly parallel non-smooth convex optimization. Wallach H, Larochelle H, Beygelzimer A, d'Alché-Buc F, Fox E, Garnett R, eds. *Advances in Neural Information Processing Systems*, vol. 32 (Curran Associates, Inc., Red Hook, NY).
- [10] Chen X, Peng B (2023) Memory-query tradeoffs for randomized convex optimization. 2023 IEEE 64th Annual Sympos. Foundations Computer Sci. (FOCS) (IEEE Computer Society, Los Alamitos, CA), 1400–1413.

- [11] Feige U, Schechtman G (2002) On the optimality of the random hyperplane rounding technique for max cut. Random Structures Algorithms 20(3):403–440.
- [12] Garg S, Raz R, Tal A (2018) Extractor-based time-space lower bounds for learning. *Proc. 50th Annual ACM SIGACT Sympos. Theory Comput.* (Association for Computing Machinery, New York), 990–1002.
- [13] Grötschel M, Lovász L, Schrijver A (2012) Geometric Algorithms and Combinatorial Optimization, vol. 2 (Springer Science & Business Media, New York).
- [14] Jiang H (2021) Minimizing convex functions with integral minimizers. Marx D, ed. *Proc. 2021 ACM-SIAM Sympos. Discrete Algorithms* (SODA) (Society for Industrial and Applied Mathematics, Philadelphia), 976–985.
- [15] Jiang AX, Leyton-Brown K (2015) Polynomial-time computation of exact correlated equilibrium in compact games. *Games Econom. Behav.* 91:347–359.
- [16] Jiang H, Lee YT, Song Z, Wong SCw (2020) An improved cutting plane method for convex optimization, convex-concave games, and its applications. Proc. 52nd Annual ACM SIGACT Sympos. Theory Comput. (Association for Computing Machinery, New York), 944–953.
- [17] Kol G, Raz R, Tal A (2017) Time-space hardness of learning sparse parities. *Proc. 49th Annual ACM SIGACT Sympos. Theory Comput.* (Association for Computing Machinery, New York), 1067–1080.
- [18] Lee YT, Sidford A, Wong SCw (2015) A faster cutting plane method and its implications for combinatorial and convex optimization. 2015 IEEE 56th Annual Sympos. Foundations Comput. Sci. (IEEE, Piscataway, NJ), 1049–1065.
- [19] Levin AY (1965) An algorithm for minimizing convex functions. Doklady Akademii Nauk SSSR 160(6):1244-1247.
- [20] Lewis AS, Overton ML (2013) Nonsmooth optimization via quasi-Newton methods. Math. Programming 141(1):135–163.
- [21] Liu DC, Nocedal J (1989) On the limited memory BFGS method for large scale optimization. Math. Programming 45(1):503-528
- [22] Marsden A, Sharan V, Sidford A, Valiant G (2022) Efficient convex optimization requires superlinear memory. Loh P-L, Raginsky M, eds. *Proc. Thirty Fifth Conf. Learn. Theory*, vol. 178 (PMLR, New York), 2390–2430.
- [23] McCormick ST (2005) Submodular function minimization. Handbooks Oper. Res. Management Sci. 12:321–391.
- [24] Mitliagkas I, Caramanis C, Jain P (2013) Memory limited, streaming PCA. Burges CJ, Bottou L, Welling M, Ghahramani Z, Weinberger KQ, eds. *Advances in Neural Information Processing Systems*, vol. 26 (Curran Associates, Inc., Red Hook, NY), 2886–2894.
- [25] Moshkovitz D, Moshkovitz M (2017) Mixing implies lower bounds for space bounded learning. Kale S, Shamir O, eds. *Proc.* 2017 Conf. Learn. Theory, vol. 65 (PMLR, New York), 1516–1566.
- [26] Moshkovitz D, Moshkovitz M (2018) Entropy samplers and strong generic lower bounds for space bounded learning. Karlin AR, ed. 9th Innovations Theoret. Comput. Sci. Conf. (ITCS 2018), Leibniz International Proceedings in Informatics (LIPIcs), vol. 94 (Schloss Dagstuhl-Leibniz-Zentrum für Informatik, Dagstuhl, Germany), 1–20.
- [27] Nemirovski A (1994) On parallel complexity of nonsmooth convex optimization. J. Complexity 10(4):451–463.
- [28] Nemirovsky A, Yudin D, Dawson E (1983) Problem Complexity and Method Efficiency in Optimization (Wiley, New York).
- [29] Nesterov JE (1989) Self-concordant functions and polynomial-time methods in convex programming. Report, Central Economic and Mathematic Institute, USSR Academy of Science, Moscow.
- [30] Nesterov Y (2003) Introductory Lectures on Convex Optimization: A Basic Course, vol. 87 (Springer Science & Business Media, New York).
- [31] Nocedal J (1980) Updating quasi-newton matrices with limited storage. Math. Comput. 35(151):773-782.
- [32] Papadimitriou CH, Roughgarden T (2008) Computing correlated equilibria in multi-player games. J. ACM 55(3):1-29.
- [33] Raz R (2017) A time-space lower bound for a large class of learning problems. 2017 IEEE 58th Annual Sympos. Foundations Comput. Sci. (FOCS) (IEEE, Piscataway, NJ), 732–742.
- [34] Sharan V, Sidford A, Valiant G (2019) Memory-sample tradeoffs for linear regression with small error. *Proc. 51st Annual ACM SIGACT Sympos. Theory Comput.* (Association for Computing Machinery, New York), 890–901.
- [35] Shor NZ (1977) Cut-off method with space extension in convex programming problems. Cybernetics 13(1):94-96.
- [36] Steinhardt J, Duchi J (2015) Minimax rates for memory-bounded sparse linear regression. Grünwald P, Hazan E, Kale S, eds. *Proc. 28th Conf. Learn. Theory*, vol. 40 (PMLR, New York), 1564–1587.
- [37] Steinhardt J, Valiant G, Wager S (2016) Memory, communication, and statistical queries. Feldman V, Rakhlin A, Shamir O, eds. 29th Annual Conf. Learn. Theory, vol. 49 (PMLR, New York), 1490–1516.
- [38] Tarasov SP, Khachiyan LG, Erlikh II (1988) The method of inscribed ellipsoids. Soviet Math. Doklady 37:226-230 [English translation].
- [39] Vaidya PM (1996) A new algorithm for minimizing convex functions over convex sets. Math. Programming 73(3):291–341.
- [40] Woodworth BE, Srebro N (2016) Tight complexity bounds for optimizing composite objectives. Lee D, Sugiyama M, Luxburg U, Guyon I, Garnett R, eds. *Advances in Neural Information Processing Systems*, vol. 29 (Curran Associates, Inc., Red Hook, NY), 3646–3654.
- [41] Woodworth BE, Srebro N (2017) Lower bound for randomized first order convex optimization. Preprint, submitted September 11, https://arxiv.org/abs/1709.03594.
- [42] Woodworth B, Srebro N (2019) Open problem: The oracle complexity of convex optimization with limited memory. Beygelzimer A, Hsu D, eds. *Proc. Thirty-Second Conf. Learn. Theory*, vol. 99 (PMLR, New York), 3202–3210.
- [43] Yudin DB, Nemirovskii AS (1976) Informational complexity and efficient methods for the solution of convex extremal problems. *Matekon* 13(2):22–45.