

Note I: Efficient Convex Optimization Requires Superlinear Memory [MSSV22]

Question: Is there a convex optimization task for which limited memory causes a need for more queries?

The answer is YES!

To show: an information-theoretic source of hardness and reduction^①
from any first-order optimization problem in a given class^②
& Exhibiting a concrete function in that class^③

The key takeaway: learning a global feature of some structure in high-dimension is difficult with little memory

For us, we focus on the null-space of a random matrix to start with^①, following "Orthogonal Vector Game"

Game 1 (OVG) Input: $d, k, B_d \subset \mathbb{R}^d, M$

Phase 1: Store Message

Oracle: For $n \leftarrow d/2$, sample $A \sim \text{Unif}(B_d^n)$

Oracle: sample random bit string of length $3 \cdot 2^d$ & present to Player
Player: observing A, R , store Message $\in \{0,1\}^M$

Phase 2: Adaptive Queries

for $i \in [M]$:

Player: based on prev. queries & responses, submits $x_i \in S_d$

Oracle: responds with predetermined $g_A(x_i) =: z_i$

Phase 3: Output

Player produces (y_1, \dots, y_n) sequence
that is approx. orthogonal to A &
robustly linearly independent.

Question: What makes the game hard?

Roughly speaking, we attempt to "find" k -different directions in $\text{null}(A)$: this should be incompressible in the following sense:

Definition (Memory-sensitive base) A sequence of sets $\{B_d\}_{d \geq 0}$, $B_d \subset \mathbb{R}^d$ is (k, c_B) -memory-sensitive if, $\forall Z = [z_1, \dots, z_k]$ projection matrices

$$\Pr[\|Z^T h\|_\infty \leq \frac{1}{2}] = \Pr\left[\max_{i=1}^k |\langle z_i, h \rangle| \leq \frac{1}{2}\right] \leq 2^{-c_B k}$$

additionally, $\Pr[\|h\|_2 > d] \leq 2^{-d}$.

So, locating k -vectors nearly orthonormal to A substantially narrows the possibilities for rows a_1^T, \dots, a_n^T to at most $2^{-c_B k} |B_d|$ options: encodable in $\log |B_d| - c_B k$ bits

Indeed, memory-sensitive bases exist, such as $\mathcal{R}_2 = \{\pm 1\}^d$.

Now, we begin efforts to showing a lower bound to AVG.
... but first, what are the trivial strategies?

① No Queries: we must solve the problem in Message immediately
we can store the subspace of $\text{null}(A)$ in $\tilde{O}(dk)$ bits
Further, Shannon source coding tells us we can't hope for much better

② No Memory: simply query $m = n$ rows of A and solve directly.

We refine our lower bound to ask the following:

Is it possible to use even slightly less than dk bits of memory & make a few adaptive queries to solve AVG?

Alas, the answer is NO!

Theorem $\exists c > 0$ s.t. $\forall k \geq \Omega(\log d)$, if $M \leq cd_k$, if the player wins the OG with probability at least $\frac{1}{2}$, then $m \geq d/6$

Proof The fact that Y is successful insinuates that it greatly reduces the entropy of A at the end of Phase 1.
[We will assume that Y is astronomical, but show it is 'WLOG' later]
We consider "reconstructing" A using Y & oracle responses:

$$\begin{aligned} I(A; Y | G, R) &= I(A; f(\text{Message}, G, R) | G, R) \\ &\leq I(A; \text{Message}, G, R | G, R) \leq H(\text{Message}) = M. \end{aligned}$$

by the bottleneck caused by passing Message after Phase 1.

However,

$$\begin{aligned} I(A; Y | G, R) &= H(A | G, R) - H(A | G, R, Y) \\ &\geq (n-m) \log |B_d| - (n-m) (\log |B_d| - c_B k) \\ &= c_B (n-m) k \end{aligned}$$

using the compression & assuming that all m oracles are distinct.

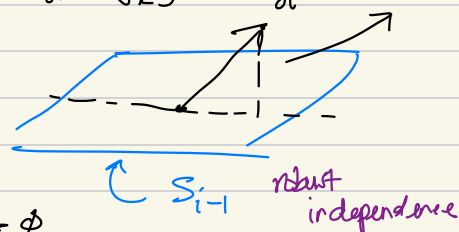
$$\text{Thus, } M \geq c_B (n-m) k \Rightarrow m \geq n - \frac{M}{c_B k}.$$

So, if $n = d/2$, $M = \frac{c_B}{3} \cdot dk \Rightarrow m \geq d/6$, as desired. \square

Note that we have oversimplified a lot of the reconstruction, but the proof is "essentially" correct.

To tie some loose ends, we formalize approximate orthogonality (a.o.) & robust linear independence (r.l.i.) and show that they do indeed imply an orthonormal basis of $\text{col}(A)$.

Definition With respect to $A \in \mathbb{R}^{d \times n}$, a sequence (y_1, \dots, y_k) is $y_i \geq \frac{1}{d}$.

$$\left[\begin{array}{l} \text{(a.o.)} \quad \frac{\|Ay_i\|_\infty}{\|y_i\|_2} \leq d^{-4} \quad \forall i, \\ \text{(r.l.i.)} \quad \text{let } S_i = \text{span}(y_1, y_2, \dots, y_i), S_0 = \emptyset, \\ \quad \quad \quad \frac{\|\text{Proj}_{S_{i-1}}(y_i)\|}{\|y_i\|_2} \leq 1 - \frac{1}{d^2} \end{array} \right.$$


Lemma: fix $\lambda \in (0, 1]$, $n_0 \leq d$. Let a sequence $(y_1, y_2, \dots, y_{n_0})$ of unit vectors satisfy λ -RLI, so

$$\|\text{Proj}_{S_{i-1}}(y_i)\| \leq 1 - \lambda.$$

Then, $\exists (m_1, m_2, \dots, m_{n_1})$ orthonormal st $\forall a \in \mathbb{R}^d$, $\|M^T a\|_\infty \leq \frac{d}{\lambda} \|\Psi^T a\|_\infty$
 Where $M = [m_1, m_2, \dots, m_{n_1}]$, $\Psi = [y_1, y_2, \dots, y_{n_0}]$ and $n_1 = \lceil n_0/2 \rceil$.

In particular, if the sequence is additionally d^{-4} -a.o. as well as d^{-2} -RLI, we have that $\|M^T a\| \leq \frac{1}{2} \leq \frac{1}{2}$, $\forall a \in \text{row}(A)$.

This completes our sketch for O/G hardness.

Next, we design an optimization task which inherently involves finding directions in the incompressible null-space of a random matrix.

We propose the following, two-part construction:

$$F(x) = \frac{1}{d^6} \max \left\{ \underbrace{d^5 \|Ax\|_\infty - 1}_{f_A(x)}, \underbrace{\max_{i=1}^N (v_i^T x - i)}_{f_v(x)} \right\}$$

Nemirovski vector.

for $A \in \{\pm 1\}^{d/2 \times d}$, $v_i \sim \frac{1}{\sqrt{d}} \{\pm 1\}^d$ drawn uniformly at random.

We define the subgradient function as

$$g_F(x) = \begin{cases} \min \left\{ i : \frac{1}{d} |a_i^T x| - \frac{1}{d^6} = F(x) \right\} & \text{if } F(x) = d^{-6} f_A(x) \\ \min \left\{ i : v_i^T x - i = d^6 F(x) \right\} & \text{else} \end{cases}$$

We also define an **informative subgradient** as a query which returns a previously unseen Nemirovski vector v_i .

Later, we will show \mathcal{F} (the distribution over (F, g_F)) has the structure of the following class.

Definition An (L, N, k, ϵ^*) -memory sensitive function class \mathcal{F} can be

sampled with at most 2^q bits, and, on any potentially randomized algorithm A_{rand} of Nd , the following holds w.p. $\geq 2/3$

1) The sampled function is L -Lipschitz

2) The informative queries are d^{-2} -r.l.i up to history of length k

2) (cont) more formally, if $\{x_{t_j}\}$ are the set of informative queries with $\{t_j\} \subset [Nd]$ and $S_j = \text{span}(\{x_{t_i} : \max(1, j-k) \leq i \leq j\})$

$$\Rightarrow \forall j \in [N], \| \text{Proj}_{S_{j-1}}(x_{t_j}) \| / \| x_{t_j} \| \leq 1 - \frac{1}{d^2}$$

3) Receiving Nemirovski vectors implies approx. orth. to A :

Any query x with $F(x) \neq \eta \|Ax\|_\infty - \rho$ sat. either

$$g_F(x) = v_i \quad \text{or} \quad \|Ax\|_\infty / \|x\|_2 \leq d^{-\eta}.$$

4) If $r < N \Rightarrow \forall i \leq t_r, F(x_i) - F(x^*) \geq \epsilon^*$

informative queries are necessary for minimization!

Immediately, let's show why this definition is useful by the reduction

Lemma Let $A \sim \text{Unif}(B_d^1)$ $(f, g_f) \sim F$ for (L, N, k, ϵ^*) -sensitive class,

If \exists an M -bit algorithm optimizing $F_{A,f}$ with at most

$M \leq N/(k+1)$ queries w.p. $\geq 2/3 \Rightarrow$ winning strategy for OG w.p. $\geq 1/2$.

Proof By Pigeonhole, there exists a period of M queries in which $k+1$ informative queries are observed.

We run the algorithm in Phase 1 by drawing $(f, g_f) \sim F$, constructing $F_{A,f}$, and saving in Message the memory state at the start of this period.

In phase 2, we continue by resampling (f, g_f) using the same randomness, and query to the oracle. We can simulate our first-order oracle again even without access to A (now?).

Finally, we search over all k -size subsets of our queries and output a robustly independent sequence.

Correctness? We can condition on properties in the memory-sensitive class with extra failure probability of $1/3$.

The existence of a solution, the number of informative queries is direct and completes the proof. \square

Let's finish by sketching the argument for showing F is memory sensitive.

The key is to show that a "resisting oracle" which draws v_j only when needed, is consistent with itself.

\Rightarrow we can show that, prior to v_j being drawn, all prior queries are independent: thus the following event

$$E = \left\{ \begin{array}{l} \forall x \text{ submitted up to phase } j, |\langle x, v_j \rangle| \leq \sqrt{\frac{10 \log d}{d}} \\ \| \text{Proj}_{S_j} v_j \| \leq \sqrt{\frac{30 k \log d}{d}} \end{array} \right.$$

has probability $\Pr[E] \geq 1 - 1/d$, by concentration

Using this, we show that the $1/\delta$ drop-off between terms is sufficient to discovering v_j 's in order, as long as $\delta \geq \sqrt{\frac{ck \log d}{d}}$.

Other properties readily follow from algebra and conditioning on E .