

# Open Problem: The Oracle Complexity of Convex Optimization with Limited Memory

**Blake Woodworth**

**Nathan Srebro**

6045 S Kenwood Ave, Chicago, IL, 60637

BLAKE@TTIC.EDU

NATI@TTIC.EDU

**Editors:** Alina Beygelzimer and Daniel Hsu

## Abstract

We note that known methods achieving the optimal oracle complexity for first order convex optimization require quadratic memory, and ask whether this is necessary, and more broadly seek to characterize the minimax number of first order queries required to optimize a convex Lipschitz function subject to a memory constraint.

## 1. Introduction

We consider first-order optimization methods for convex Lipschitz bounded functions. I.e., for the following optimization problem

$$\min_{x \in \mathbb{R}^d: \|x\| \leq B} F(x) \quad (1)$$

where  $F$  is convex and  $L$ -Lipschitz, we consider using (exact) queries returning  $F(x)$  and a subgradient  $\nabla F(x)$ , and ask the classical question of how many queries are required to ensure we find an  $\epsilon$ -suboptimal solution. The classical answer is that  $O(d \log LB/\epsilon)$  suffice, using the center-of-mass method, and that, when  $d \ll (LB/\epsilon)^2$  this is optimal (Nemirovsky and Yudin, 1983). However, the center-of-mass method is intractable, at least exactly. Other methods with a  $O(\text{poly}(d) \text{polylog}(LB/\epsilon))$ , and even  $\tilde{O}(d \log(LB/\epsilon))$ , query complexity and polynomial runtime have been suggested, including as the Ellipsoid method (Shor, 1970), Vaidya’s method (Atkinson and Vaidya, 1995) and approximate center of mass using sampling (Bertsimas and Vempala, 2002). But these are generally not used in practice, since the higher order polynomial runtime dependence on the dimension is prohibitive. These methods also all require storing all returned gradients, or alternatively an ellipsoid in  $\mathbb{R}^d$ , and so  $\Omega(d^2)$  memory. A simpler alternative is gradient descent, which requires  $O((LB/\epsilon)^2)$  queries, but only  $O(d \log LB/\epsilon)$  memory and  $O(d)$  runtime per query.

One might ask: is it possible to achieve the optimal query complexity using a “simple” method? Since it is much harder to provide runtime lower bound, we instead focus on the required memory, and ask: **is it possible to achieve the optimal query complexity with  $O(d \log LB/\epsilon)$  memory? How does the first-order oracle complexity trade off with the memory needed by an optimization algorithm?** This question is formalized in the following Section.

## 2. Problem Formulation

We capture the class of first-order optimization algorithms that use  $M$  bits of memory in terms of a set of “encoders” and “decoders.” In each iteration, the decoder reads the  $M$  bits of memory, and determines a query point  $x_t$ . The encoder receives the function value  $F(x_t)$  and a subgradient  $\nabla F(x_t)$  at  $x_t$ —as is standard with oracle based optimization, if  $F(x_t)$  is not differentiable at  $x_t$ , we require the method works for any valid subgradient used. The encoder then uses the current memory state,  $F(x_t)$  and  $\nabla F(x_t)$  to update the memory state for the next iteration. At the end, the algorithm’s output is chosen as a function of the final memory state. To be clear, the encoding and decoding functions can require an arbitrary amount of memory to compute, and can compute using real numbers. However, between each access to the oracle, there is a “bottleneck” where the algorithm’s state must be compressed down to  $M$  bits.

Formally, we define  $\mathcal{A}_{M,T}$ , the class of all deterministic<sup>1</sup> first-order optimization algorithms that use  $M$  bits of memory and  $T$  function value and gradient computations. An algorithm  $A \in \mathcal{A}_{M,T}$  is specified by a set of decoder functions  $\{\phi_t : \{0, 1\}^M \rightarrow \mathbb{R}^d\}_{t=1}^T$ , a set of encoder functions  $\{\psi_t : \{0, 1\}^M \times \mathbb{R} \times \mathbb{R}^d \rightarrow \{0, 1\}^M\}_{t=1}^T$ , and an output function  $\zeta : \{0, 1\}^M \rightarrow \mathbb{R}^d$ . The algorithm’s memory is initially blank  $\mu_1 = 0$ . The  $T$  iteration  $t = 1 \dots T$  are specified recursively by  $x_t = \phi_t(\mu_t)$  and

$$\mu_{t+1} = \psi_t(\mu_t, F(x_t), \nabla F(x_t)) = \psi_t(\mu_t, F(\phi_t(\mu_t)), \nabla F(\phi_t(\mu_t))), \quad (2)$$

and the output of the algorithm, denoted  $F(A)$ , is given by:

$$A(F) = \zeta(\mu_T). \quad (3)$$

Let  $\mathcal{F}_{L,B}^d$  be the set of all convex,  $L$ -Lipschitz functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\exists x^* \in \arg \min_x F(x)$  with  $\|x^*\| \leq B$ . We define the minimax memory-bounded first-order oracle complexity as

$$T_{L,B}(d, M, \epsilon) = \inf \left\{ T \in \mathbb{N} : \inf_{A \in \mathcal{A}_{M,T}} \sup_{F \in \mathcal{F}_{L,B}^d} F(A(F)) - \min_{\|x\| \leq B} F(x) \leq \epsilon \right\}, \quad (4)$$

where the supremum over functions  $F$  should be interpreted also as a supremum over all valid subgradients  $\nabla F(x_t)$  used in the updates. Without loss of generality, we will fix  $L = B = 1$  and write  $T(d, M, \epsilon) := T_{1,1}(d, M, \epsilon)$ . We will further say that a query-memory tradeoff  $(T, M)$  is *possible* for a problem specified by  $(d, \epsilon)$  if  $T \geq T(d, M, \epsilon)$  and *impossible* if  $T < T(d, M, \epsilon)$ .

## 3. Current Knowledge

In high dimensions, when  $d = \Omega(\frac{1}{\epsilon^2 \log 1/\epsilon})$ , gradient descent is optimal in terms of both query and memory complexity, and so we consider only  $d = O(1/\epsilon^2)$ .

1. We focus on deterministic algorithms, but an analogous class of randomized algorithms could easily be specified. However, it seems unnecessary to complicate things in this way because there is evidence that there is little to be gained through randomization for solving problems of the form (1) [Woodworth and Srebro \(2017\)](#).

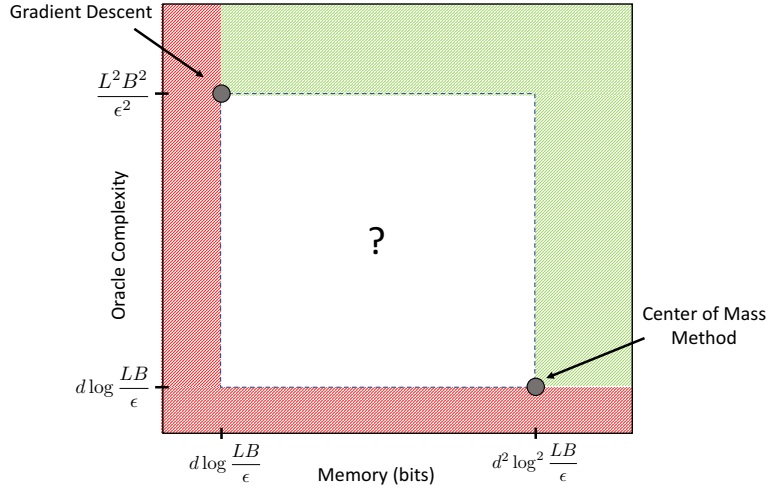


Figure 1: The tradeoff between the memory needed and the first-order oracle (query) complexity of optimization algorithms. The shaded red “L” shaped region along the bottom and left are trade-offs we know are impossible. The shaded green inverted-“L” shaped region along the top and right are trade-offs we know are possible. We do not know whether any trade-offs inside the “?” square are possible or not.

We can describe the minimax complexity  $T(d, M, \epsilon)$ , and the query-memory tradeoff, in terms of the regions of possible and impossible  $(T, M)$ , as depicted in Figure 2. We currently understand only the extremes. With any amount of memory,  $T = \Omega(d \log 1/\epsilon)$  queries are required, providing a lower bound for the possible region in terms of the query complexity (a horizontal lower bound in Figure 2). This is attained by the center of mass method, using  $O(d^2 \log^2(1/\epsilon))$  bits of memory (see Appendix B for an analysis of Center of Mass with discrete memory), and so any  $(T = \Omega(d \log 1/\epsilon), M = \Omega(d^2 \log^2 1/\epsilon))$  is possible (the rectangle above and to the right of “Center of Mass” in Figure 2). At the other extreme, even just representing the answer requires  $\Omega(d \log(LB/\epsilon))$  bits (see Theorem 5 in Appendix C), providing a lower bound for the possible region in terms of memory (the vertical lower bound in Figure 2). This is attained by Gradient Descent using  $O(1/\epsilon^2)$  queries (see Appendix A for an analysis of Gradient Descent with discrete memory), and so any  $(T = \Omega(1/\epsilon^2), M = \Omega(d \log 1/\epsilon))$  is possible (the rectangle above and to the left of “Gradient Descent” in Figure 2).

To the best of our knowledge, what happens inside the square bordered by these regions is completely unknown. Nothing we know would contradict the existence of a  $T, M = O(d \log 1/\epsilon)$  query and memory complexity algorithm, i.e. a single optimal method at the bottom left corner of the unknown square, making the entire square possible. It is also entirely possible, as far as we know, that improving over a query complexity of  $\Theta(1/\epsilon^2)$  requires  $\Omega(d^2 \log 1/\epsilon)$  memory, making the entire square impossible, and implying that no compromise is possible between the query requirement of Gradient Descent and memory requirement of Center of Mass.

## 4. Challenges

Ultimately, we would like fully understand what is and is not possible:

**Question 1 (\$500 or Two Star Michelin Meal)** Provide a complete characterization of  $T(d, M, \epsilon)$  and the possible  $(T, M)$  trade-off, preferably up to constant factors, and at most up to factors poly-logarithmic in  $T$  and  $M$ .

The most interesting scaling of  $d$  and  $\epsilon$  is when the dimension is larger than poly-logarithmic but smaller than polynomial in  $1/\epsilon$ , so that  $d \log 1/\epsilon$  memory is less than quadratic memory, but  $1/\epsilon^2$  query complexity is not polynomial in  $d$ .

Even without understanding the entire trade-off, it would be interesting to study what can be done on its boundary. Perhaps the most important regime is the case of linear memory  $M = \Theta(d \log \frac{LB}{\epsilon})$ . Therefore, as a starting point, we ask to characterize  $T(d, \Theta(d \text{polylog} \frac{LB}{\epsilon}), 1/\epsilon)$ . In particular, is it possible to have query complexity polynomial in  $d$  with  $\tilde{O}(d)$  memory?

**Question 2 (\$200 or One Star Michelin Meal)** Can we have  $T(d, M = \tilde{O}(d), \epsilon) = O(\text{poly } d)$  when  $d = \Omega(\log^c 1/\epsilon)$  but  $d = O(1/\epsilon^c)$  for all  $c$ ?

At the other extreme, we might ask whether quadratic memory is necessary in order to achieve optimal query complexity:

**Question 3 (\$200 or One Star Michelin Meal)** Can we have  $T(d, M = O(d^{2-\delta}), \epsilon) = \tilde{O}(d \text{polylog } 1/\epsilon)$ , for  $\delta > 0$ , when  $d = \Omega(\log^c 1/\epsilon)$  but  $d = O(1/\epsilon^c)$  for all  $c$ ?

The above represent specific incursions into the unknown square in Figure 2. Any other such incursion would also be interesting, and provide either for a memory lower bound, or a trade-off improving over Gradient Descent and Center of Mass in some regime.

**Question 4 (\$100 or Michelin Bib Gourmand Meal)** Resolve the possibility or impossibility of some trade-off  $(T, M)$  polynomially inside the unknown square in Figure 2.

## References

- David S Atkinson and Pravin M Vaidya. A cutting plane algorithm for convex programming that uses analytic centers. *Mathematical Programming*, 69(1-3):1–43, 1995.
- Dimitris Bertsimas and Santosh Vempala. Solving convex programs by random walks. In *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*, pages 109–115. ACM, 2002.
- Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- Branko Grünbaum et al. Partitions of mass-distributions and of convex bodies by hyperplanes. *Pacific Journal of Mathematics*, 10(4):1257–1261, 1960.

Arkadii Semenovich Nemirovsky and David Borisovich Yudin. Problem complexity and method efficiency in optimization. 1983.

Naum Z Shor. Convergence rate of the gradient descent method with dilatation of the space. *Cybernetics*, 6(2):102–108, 1970.

Blake Woodworth and Nathan Srebro. Lower bound for randomized first order convex optimization. *arXiv preprint arXiv:1709.03594*, 2017.

## Appendix A. Analysis of Gradient Descent

---

### Algorithm 1 Memory-Bounded Gradient Descent

---

Initialize:  $x_0 = 0, x_{\text{best}} = 0, F_{\text{best}} = \infty, \mu_0 = \{x_0, x_{\text{best}}, F_{\text{best}}\}$   
**for**  $t = 0, \dots, T$  **do**  
    Compute  $F(x_t), \nabla F(x_t)$  }  $\phi_t(\mu_t)$   
    **if**  $F(x_t) < F_{\text{best}}$  **then**  
         $F_{\text{best}} = \text{Discretize}(F(x_t))$  }  $\psi_t(\mu_t, \phi(\mu_t))$   
         $x_{\text{best}} = x_t$  }  $\psi_t(\mu_t, \phi(\mu_t))$   
         $\tilde{x}_{t+1} = x_t - \eta_t \nabla F(x_t)$   
         $x_{t+1} = \text{Discretize}(\tilde{x}_{t+1})$   
         $\mu_{t+1} = \{x_{t+1}, x_{\text{best}}, F_{\text{best}}\}$   
**return**  $x_{\text{best}}$

---

**Theorem 1** For any  $L$ -Lipschitz and convex function  $F$  with  $\|x^*\| \leq B$ , the gradient descent algorithm can find a point  $\hat{x}$  with  $F(\hat{x}) - F^* \leq \epsilon$  using  $O(d \log \frac{LB}{\epsilon})$  bits of memory and  $O\left(\frac{L^2 B^2}{\epsilon^2}\right)$  function and gradient evaluations.

### Proof

For now, assume that in each iteration the perturbation of the gradient descent iterates resulting from the discretization  $\tilde{x}_t \mapsto x_t$  is bounded in L2 norm, i.e.  $\|\tilde{x}_t - x_t\| \leq D$ . Then, following the standard gradient descent analysis,

$$\|x_{t+1} - x^*\|^2 = \|\tilde{x}_{t+1} - x^* + x_{t+1} - \tilde{x}_{t+1}\|^2 \tag{5}$$

$$\leq \|x_t - \eta_t \nabla F(x_t) - x^*\|^2 + D^2 + 2 \langle \tilde{x}_{t+1} - x^*, x_{t+1} - \tilde{x}_{t+1} \rangle \tag{6}$$

$$= \|x_t - x^*\|^2 + \eta_t^2 \|\nabla F(x_t)\|^2 - 2\eta_t \langle \nabla F(x_t), x_t - x^* \rangle + D^2 + 2DB \tag{7}$$

$$\leq \|x_t - x^*\|^2 + \eta_t^2 L^2 - 2\eta_t (F(x_t) - F^*) + D^2 + 2DB \tag{8}$$

Rearranging this expression, we conclude

$$F(x_t) - F^* \leq \frac{1}{2\eta_t} \left( \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 \right) + \frac{\eta_t L^2}{2} + \frac{D^2 + 2DB}{2\eta_t} \tag{9}$$

Choosing a fixed stepsize  $\eta_t = \eta = \frac{B}{L\sqrt{T}}$  and averaging the iterates, we conclude  $\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t$  achieves suboptimality

$$F(\bar{x}_T) - F^* \leq \frac{1}{T} \sum_{t=1}^T F(x_t) - F^* \quad (10)$$

$$\leq \frac{1}{T} \sum_{t=1}^T \frac{1}{2\eta} \left( \|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2 \right) + \frac{\eta L^2}{2} + \frac{D^2 + 2DB}{2\eta} \quad (11)$$

$$= \frac{L}{2B\sqrt{T}} \left( \|x_1 - x^*\|^2 - \|x_{T+1} - x^*\|^2 \right) + \frac{LB}{2\sqrt{T}} + \frac{(D^2 + 2DB)L\sqrt{T}}{2B} \quad (12)$$

$$\leq \frac{LB}{\sqrt{T}} + \frac{(D^2 + 2DB)L\sqrt{T}}{2B} \quad (13)$$

Thus,  $D \leq \frac{B}{T}$  ensures

$$F(\bar{x}_T) - F^* \leq \frac{LB}{\sqrt{T}} + \frac{(D^2 + 2DB)L\sqrt{T}}{2B} \quad (14)$$

$$\leq \frac{LB}{\sqrt{T}} + \left( \frac{B}{2T^2} + \frac{2B}{2T} \right) L\sqrt{T} \quad (15)$$

$$\leq \frac{3LB}{\sqrt{T}} \quad (16)$$

Since the averaged iterate achieves this suboptimality, the best iterate's suboptimality is at least this good. For  $T \geq \frac{9L^2B^2}{\epsilon^2}$ , this ensures that at least one of the iterates was  $\epsilon$ -suboptimal. As long as  $F_{\text{best}}$  is discretized to accuracy  $\epsilon$ , then  $x_{\text{best}}$  is at most  $2\epsilon$ -suboptimal. This discretization can be achieved using  $\log \frac{2LB}{\epsilon}$  bits.

Discretizing the iterates up to accuracy  $D = \frac{B}{T} = \frac{\epsilon^2}{9L^2B}$  can be achieved using the log of the  $\frac{\epsilon^2}{9L^2B}$  L2 covering number of the radius- $B$  ball, which is upper bounded by  $d \log \left( 1 + \frac{36L^2B^2}{\epsilon^2} \right)$  bits. The discretization of  $x_{\text{best}}$  is achieved using the same number of bits.

Therefore, the total number of bits of memory needed to implement gradient descent is at most

$$2 \cdot d \log \left( 1 + \frac{36L^2B^2}{\epsilon^2} \right) + \log \frac{2LB}{\epsilon} = O \left( d \log \frac{LB}{\epsilon} \right) \quad (17)$$

■

## Appendix B. Analysis of Center of Mass Algorithm

**Lemma 2 (Grünbaum et al. (1960))** *For any convex set  $K \subseteq \mathbb{R}^d$  with center of gravity  $c$ , and any halfspace  $H = \{x : \langle a, x - c \rangle \geq 0\}$  passing through  $c$ ,*

$$\frac{1}{e} \leq \frac{\text{Vol}(K \cap H)}{\text{Vol}(K)} \leq 1 - \frac{1}{e}$$

---

**Algorithm 2** Memory-Bounded Center of Mass
 

---

 Initialize:  $K_0 = \{x : \|x\| \leq B\}$ ,  $\mu_0 = \emptyset$ 
**for**  $t = 0, \dots, T$  **do**
**for**  $k = 1, \dots, t$  **do**

$$\left. \begin{aligned}
 c_{k-1} &= \int_{K_{k-1}} x dx / \int_{K_{k-1}} dx \\
 K_k &= K_{k-1} \cap \left\{ x : \left\langle \tilde{\nabla} F(c_{k-1}), x - c_{k-1} \right\rangle \leq 0 \right\} \\
 c_t &= \int_{K_t} x dx / \int_{K_t} dx \\
 \tilde{F}(c_t), \tilde{\nabla} F(c_t) &= \text{Discretize}(F(c_t), \nabla F(c_t)) \\
 \mu_{t+1} &= \mu_t \cup \left\{ \tilde{F}(c_t), \tilde{\nabla} F(c_t) \right\}
 \end{aligned} \right\} \begin{array}{l} \phi_t(\mu_t) \\ \psi_t(\mu_t, \phi(\mu_t)) \end{array}$$

**return**  $\text{Discretize}\left(\arg \min_{c \in \{c_1, \dots, c_T\}} \tilde{F}(c)\right)$ 


---

**Theorem 3** For any  $L$ -Lipschitz and convex function  $F$  with  $\|x^*\| \leq B$ , the center of mass algorithm can find a point  $\hat{x}$  with  $F(\hat{x}) - F^* \leq \epsilon$  using  $O(d^2 \log^2 \frac{LB}{\epsilon})$  bits of memory and  $O(d \log \frac{LB}{\epsilon})$  function and gradient evaluations.

**Proof** This proof is quite similar to existing analysis of the center of mass algorithm (Bubeck et al., 2015), we simply take care to count the number of required bits.

Consider the set  $K^\alpha = \{(1 - \alpha)x^* + \alpha x : x \in K_0\}$ , which has volume  $\text{Vol}(K^\alpha) = \alpha^d \text{Vol}(K_0)$ . By convexity,

$$F((1 - \alpha)x^* + \alpha x) \leq (1 - \alpha)F^* + \alpha F(x) \quad (18)$$

$$\leq (1 - \alpha)F^* + \alpha(F^* + \|\nabla F(x)\| \|x - x^*\|) \quad (19)$$

$$\leq F^* + 2\alpha LB \quad (20)$$

By Grünbaum's Lemma,  $\text{Vol}(K_T) \leq (1 - \frac{1}{e}) \text{Vol}(K_{T-1}) \leq (1 - \frac{1}{e})^T \text{Vol}(K_0)$ . Thus, when  $T \geq 3d \log(1/\alpha)$

$$\text{Vol}(K_T) \leq \left(1 - \frac{1}{e}\right)^T \text{Vol}(K_0) < \alpha^d \text{Vol}(K_0) = \text{Vol}(K^\alpha) \quad (21)$$

We conclude that there must be some iteration  $t$  in which  $\exists y \in K^\alpha \cap (K_t \setminus K_{t+1})$ . Thus,  $y \in K^\alpha$  and  $\left\langle \tilde{\nabla} F(c_t), y - c_t \right\rangle > 0$ . We will now argue that the center of mass  $c_t$  has small error:

$$F(c_t) \leq F(y) + \langle \nabla F(c_t), c_t - y \rangle \quad (22)$$

$$\leq F^* + 2\alpha LB + \left\langle \nabla F(c_t) - \tilde{\nabla} F(c_t), c_t - y \right\rangle + \left\langle \tilde{\nabla} F(c_t), c_t - y \right\rangle \quad (23)$$

$$\leq F^* + 2\alpha LB + \left\| \nabla F(c_t) - \tilde{\nabla} F(c_t) \right\| \|c_t - y\| + 0 \quad (24)$$

$$\leq F^* + 2B \left( \alpha L + \left\| \nabla F(c_t) - \tilde{\nabla} F(c_t) \right\| \right) \quad (25)$$

Therefore, if we choose  $\alpha = \frac{\epsilon}{4LB}$  and discretize gradients with L2 error at most  $\frac{\epsilon}{4B}$ , this ensures that  $F(c_t) - F^* \leq \epsilon$ .

The gradients of an  $L$ -Lipschitz function are contained in the Euclidean ball of radius  $L$ . Therefore, the gradients can be discretized with error  $\frac{\epsilon}{4B}$  using the logarithm of the  $\frac{\epsilon}{4B}$  L2 covering number of the Euclidean ball of radius  $L$ , which is upper bounded by  $d \log \left(1 + \frac{16LB}{\epsilon}\right)$  bits. There are  $T = 3d \log(1/\alpha) = 3d \log \left(\frac{4LB}{\epsilon}\right)$  gradients in total, thus the total number of bits required to represent the gradients is

$$T \cdot d \log \left(1 + \frac{16LB}{\epsilon}\right) = 3d^2 \log \left(\frac{4LB}{\epsilon}\right) \log \left(1 + \frac{16LB}{\epsilon}\right) \quad (26)$$

Since  $\epsilon \leq LB$ , this is upper bounded by  $3d^2 \log^2 \left(\frac{17LB}{\epsilon}\right)$  bits.

Once  $T$  iterations have been completed, we know that at least one of the centers must be an  $\epsilon$ -approximate minimizer of the objective. Using the stored gradients, we can then recompute all centers and return a discretization of the best center. As long as  $\left|\tilde{F}(c_t) - F(c_t)\right| \leq \epsilon$  for all  $t$ , then the center chosen by the algorithm will be within  $\epsilon$  of the best center. This discretization of the function values requires  $T \cdot \log \left(\frac{2LB}{\epsilon}\right)$  bits.

As long as the discretization of the chosen center has L2 error at most  $\epsilon/L$ , then the output will be a  $3\epsilon$ -approximate minimizer. The number of bits needed for this discretization is at most  $d \log \left(1 + \frac{4LB}{\epsilon}\right)$ . Therefore, the total number of bits needed is at most

$$3d^2 \log^2 \left(\frac{17LB}{\epsilon}\right) + 3d \log \left(\frac{4LB}{\epsilon}\right) \log \left(\frac{2LB}{\epsilon}\right) + d \log \left(1 + \frac{4LB}{\epsilon}\right) = O \left(d^2 \log^2 \left(\frac{LB}{\epsilon}\right)\right) \quad (27)$$

Rescaling  $\epsilon' = \epsilon/3$  completes the proof.  $\blacksquare$

## Appendix C. Memory Lower Bound

**Lemma 4** *The packing number of the Euclidean unit sphere in  $\mathbb{R}^d$  with distance  $\alpha$  is at least  $\alpha^{-d}$ .*

**Proof** Let  $\{x_1, \dots, x_N\}$  be the largest possible packing of the unit sphere, with  $N < \alpha^{-d}$ . Consider the set of points that are within  $\alpha$  of one of the points in the packing:

$$\text{Vol} \left( x : \exists i \|x - x_i\| \leq \alpha \right) = \text{Vol} \left( \bigcup_i B(x_i, \alpha) \right) \quad (28)$$

$$\leq \sum_i \text{Vol} (B(x_i, \alpha)) \quad (29)$$

$$= N \alpha^d \text{Vol} (B(0, 1)) \quad (30)$$

$$< \text{Vol} (B(0, 1)) \quad (31)$$

Therefore, there exists a point  $y \in B(0, 1)$  such that  $\|y - x_i\| > \alpha$  for all  $i$ . The existence of such a point contradicts the assumption that  $\{x_1, \dots, x_N\}$  is the largest possible packing. We conclude that the packing number is at least  $\alpha^{-d}$ .  $\blacksquare$



**Theorem 5** For any  $L, B > 0$  and any  $\epsilon \leq \frac{LB}{2}$ , any optimization algorithm that is guaranteed to return an  $\epsilon$ -suboptimal point for any convex,  $L$ -Lipschitz function with  $\|x^*\| \leq B$  must use at least  $d \log \frac{LB}{2\epsilon}$  bits of memory.

**Proof** To begin, by Lemma 4 there exists a packing  $\{x_1, \dots, x_N\}$  of the ball  $\{x : \|x\| \leq B\}$  of size at least  $N \geq \left(\frac{LB}{2\epsilon}\right)^d$  such that  $\|x_i - x_j\| > \frac{2\epsilon}{L}$  for all  $i \neq j \in [N]$ . We will associate a function with each point in the packing, let

$$f_i(x) = L \|x - x_i\| \tag{32}$$

These functions are convex and  $L$ -Lipschitz, and their optimizers  $x_i$  have norm less than  $B$ .

Note that any point  $x$  which is an approximate minimizer of some  $f_i$  must have high function value on all other functions  $f_j$ . Suppose  $f_i(x) \leq \epsilon$ , then  $\|x - x_i\| \leq \frac{\epsilon}{L}$ . Consequently, for all  $j \neq i$ ,  $\|x - x_j\| = \|x - x_i + x_i - x_j\| \geq \|x_i - x_j\| - \|x - x_i\| > \frac{\epsilon}{L}$ , thus  $f_j(x) = L \|x - x_j\| > \epsilon$ .

Consider using a memory-bounded optimization algorithm to optimize one of these functions  $f_i$ . After the algorithm has made all of its first-order oracle accesses, the output function  $\zeta$  must map from the final memory state  $\mu_T$  to a solution  $\hat{x}$ . Suppose the final memory state  $\mu_T$  uses  $M < \log N$  bits, then there are at most  $2^M < N$  outputs that the algorithm might give. However, as we just argued, there exist  $N$  functions such that returning an accurate solution for any one of them requires returning an inaccurate solution for all the others. Consequently, any algorithm which can output fewer than  $N$  different outputs will fail to optimize at least one of the functions  $f_1, \dots, f_N$ . ■