

Sublinear Optimization for Machine Learning

Kenneth L. Clarkson* Elad Hazan† David P. Woodruff‡

November 26, 2024

Abstract

We give sublinear-time approximation algorithms for some optimization problems arising in machine learning, such as training linear classifiers and finding minimum enclosing balls. Our algorithms can be extended to some kernelized versions of these problems, such as SVDD, hard margin SVM, and L_2 -SVM, for which sublinear-time algorithms were not known before. These new algorithms use a combination of a novel sampling techniques and a new multiplicative update algorithm. We give lower bounds which show the running times of many of our algorithms to be nearly best possible in the unit-cost RAM model. We also give implementations of our algorithms in the semi-streaming setting, obtaining the first low pass polylogarithmic space *and* sublinear time algorithms achieving arbitrary approximation factor.

1 Introduction

Linear classification is a fundamental problem of machine learning, in which positive and negative examples of a concept are represented in Euclidean space by their feature vectors, and we seek to find a hyperplane separating the two classes of vectors.

The Perceptron Algorithm for linear classification is one of the oldest algorithms studied in machine learning [Nov62, MP88]. It can be used to efficiently give a good approximate solution, if one exists, and has nice noise-stability properties which allow it to be used as a subroutine in many applications such as learning with noise [Byl94, BFKV98], boosting [Ser99] and more general optimization [DV04]. In addition, it is extremely simple to implement: the algorithm starts with an arbitrary hyperplane, and iteratively finds a vector on which it errs, and moves in the direction of this vector by adding a multiple of it to the normal vector to the current hyperplane.

The standard implementation of the Perceptron Algorithm must iteratively find a “bad vector” which is classified incorrectly, that is, for which the inner product with the current normal vector has an incorrect sign. Our new algorithm is similar to the Perceptron Algorithm, in that it maintains a hyperplane and modifies it iteratively, according to the examples seen. However, instead of explicitly finding a bad vector, we run another *dual* learning algorithm to learn the “most adversarial” distribution over the vectors, and use that distribution to generate an “expected bad” vector. Moreover, we do not compute the inner products with the current normal vector exactly, but instead estimate them using a fast sampling-based scheme.

Thus our update to the hyperplane uses a vector whose “badness” is determined quickly, but very crudely. We show that despite this, an approximate solution is still obtained in about the

*IBM Almaden Research Center, San Jose, CA

†Department of Industrial Engineering, Technion - Israel Institute of technology, Haifa 32000 Israel. Work done while at IBM Almaden Research Center

‡IBM Almaden Research Center, San Jose, CA

Problem	Previous time	Time Here	Lower Bound
classification/perceptron	$\tilde{O}(\varepsilon^{-2}M)$ [Nov62]	$\tilde{O}(\varepsilon^{-2}(n+d))$ §2	$\Omega(\varepsilon^{-2}(n+d))$ §7.1
min. enc. ball (MEB)	$\tilde{O}(\varepsilon^{-1/2}M)$ [SV09]	$\tilde{O}(\varepsilon^{-2}n + \varepsilon^{-1}d)$ §3.1	$\Omega(\varepsilon^{-2}n + \varepsilon^{-1}d)$ §7.2
QP in the simplex	$O(\varepsilon^{-1}M)$ [FW56]	$\tilde{O}(\varepsilon^{-2}n + \varepsilon^{-1}d)$ §3.3	
Las Vegas versions		additive $O(M)$ Cor 2.11	$\Omega(M)$ §7.4
kernelized MEB and QP		factors $O(s^4)$ or $O(q)$ §6	

Figure 1: Our results, except for semi-streaming and parallel

same number of iterations as the standard perceptron. So our algorithm is faster; notably, it can be executed in time *sublinear* in the size of the input data, and still have good output, with high probability. (Here we must make some reasonable assumptions about the way in which the data is stored, as discussed below.)

This technique applies more generally than to the perceptron: we also obtain sublinear time approximation algorithms for the related problems of finding an approximate Minimum Enclosing Ball (MEB) of a set of points, and training a Support Vector Machine (SVM), in the hard margin or L_2 -SVM formulations.

We give lower bounds that imply that our algorithms for classification are best possible, up to polylogarithmic factors, in the unit-cost RAM model, while our bounds for MEB are best possible up to an $\tilde{O}(\varepsilon^{-1})$ factor. For most of these bounds, we give a family of inputs such that a single coordinate, randomly “planted” over a large collection of input vector coordinates, determines the output to such a degree that all coordinates in the collection must be examined for even a $2/3$ probability of success.

We show that our algorithms can be implemented in the parallel setting, and in the semi-streaming setting; for the latter, we need a careful analysis of arithmetic precision requirements and an implementation of our primal-dual algorithms using lazy updates, as well as some recent sampling technology [MW10].

Our approach can be extended to give algorithms for the kernelized versions of these problems, for some popular kernels including the Gaussian and polynomial, and also easily gives Las Vegas results, where the output guarantees always hold, and only the running time is probabilistic.¹ Our approach also applies to the case of soft margin SVM (joint work in progress with Nati Srebro).

Our main results, except for semi-streaming and parallel algorithms, are given in Figure 1. The notation is as follows. All the problems we consider have an $n \times d$ matrix A as input, with M nonzero entries, and with each row of A with Euclidean length no more than one. The parameter $\varepsilon > 0$ is the additive error; for MEB, this can be a relative error, after a simple $O(M)$ preprocessing step. We use the asymptotic notation $\tilde{O}(f) = O(f \cdot \text{polylog} \frac{nd}{\varepsilon})$. The parameter σ is the *margin* of the problem instance, explained below. The parameters s and q determine the standard deviation of a Gaussian kernel, and degree of a polynomial kernel, respectively.

The time bounds given for our algorithms, except the Las Vegas ones, are under the assumption of constant error probability; for output guarantees that hold with probability $1 - \delta$, our bounds should be multiplied by $\log(n/\delta)$.

The time bounds also require the assumption that the input data is stored in such a way that a given entry $A_{i,j}$ can be recovered in constant time. This can be done by, for example, keeping

¹For MEB and the kernelized versions, we assume that the Euclidean norms of the relevant input vectors are known. Even with the addition of this linear-time step, all our algorithms improve on prior bounds, with the exception of MEB when $M = o(\varepsilon^{-3/2}(n+d))$.

each row A_i of A as a hash table. (Simply keeping the entries of the row in sorted order by column number is also sufficient, incurring an $O(\log d)$ overhead in running time for binary search.)

By appropriately modifying our algorithms, we obtain algorithms with very low pass, space, and time complexity. Many problems cannot be well-approximated in one pass, so a model permitting a small number of passes over the data, called the semi-streaming model, has gained recent attention [FKM⁺08, Mut05]. In this model the data is explicitly stored, and the few passes over it result in low I/O overhead. It is quite suitable for problems such as MEB, for which any algorithm using a single pass and sublinear (in n) space cannot approximate the optimum value to within better than a fixed constant [AS10]. Unlike traditional semi-streaming algorithms, we also want our algorithms to be sublinear time, so that in each pass only a small portion of the input is read.

We assume we see the points (input rows) one at a time in an arbitrary order. The space is measured in bits. For MEB, we obtain an algorithm with $\tilde{O}(\varepsilon^{-1})$ passes, $\tilde{O}(\varepsilon^{-2})$ space, and $\tilde{O}(\varepsilon^{-3}(n+d))$ total time. For linear classification, we obtain an algorithm with $\tilde{O}(\varepsilon^{-2})$ passes, $\tilde{O}(\varepsilon^{-2})$ space, and $\tilde{O}(\varepsilon^{-4}(n+d))$ total time. For comparison, prior streaming algorithms for these problems [AS10, ZZC06] require a prohibitive $\Omega(d)$ space, and none achieved a sublinear $o(nd)$ amount of time. Further, their guarantee is an approximation up to a fixed constant, rather than for a general ε (though they can achieve a single pass).

Formal Description: Classification In the linear classification problem, the learner is given a set of n labeled examples in the form of d -dimensional vectors, comprising the input matrix A . The labels comprise a vector $y \in \{+1, -1\}^n$.

The goal is to find a separating hyperplane, that is, a normal vector x in the unit Euclidean ball \mathbb{B} such that for all i , $y(i) \cdot A_i x \geq 0$; here $y(i)$ denotes the i 'th coordinate of y . As mentioned, we will assume throughout that $A_i \in \mathbb{B}$ for all $i \in [n]$, where generally $[m]$ denotes the set of integers $\{1, 2, \dots, m\}$.

As is standard, we may assume that the labels $y(i)$ are all 1, by taking $A_i \leftarrow -A_i$ for any i with $y(i) = -1$. The approximation version of linear classification (which is necessary in case there is noise), is to find a vector $x_\varepsilon \in \mathbb{B}$ that is an ε -approximate solution, that is,

$$\forall i' \quad A_{i'} x_\varepsilon \geq \max_{x \in \mathbb{B}} \min_i A_i x - \varepsilon. \quad (1)$$

The optimum for this formulation is obtained when $\|x\| = 1$, except when no separating hyperplane exists, and then the optimum x is the zero vector.

Note that $\min_i A_i x = \min_{p \in \Delta} p^\top A x$, where $\Delta \subset \mathbb{R}^n$ is the unit simplex $\{p \in \mathbb{R}^n \mid p_i \geq 0, \sum_i p_i = 1\}$. Thus we can regard the optimum as the outcome of a game to determine $p^\top A x$, between a minimizer choosing $p \in \Delta$, and a maximizer choosing $x \in \mathbb{B}$, yielding

$$\sigma \equiv \max_{x \in \mathbb{B}} \min_{p \in \Delta} p^\top A x,$$

where this optimum σ is called the *margin*. From standard duality results, σ is also the optimum of the dual problem

$$\min_{p \in \Delta} \max_{x \in \mathbb{B}} p^\top A x,$$

and the optimum vectors p^* and x^* are the same for both problems.

The classical Perceptron Algorithm returns an ε -approximate solution to this problem in $\frac{1}{\varepsilon^2}$ iterations, and total time $O(\varepsilon^{-2}M)$.

For given $\delta \in (0, 1)$, our new algorithm takes $O(\varepsilon^{-2}(n+d)(\log n) \log(n/\delta))$ time to return an ε -approximate solution with probability at least $1 - \delta$. Further, we show this is optimal in the unit-cost RAM model, up to poly-logarithmic factors.

Formal Description: Minimum Enclosing Ball (MEB) The MEB problem is to find the smallest Euclidean ball in \mathbb{R}^d containing the rows of A . It is a special case of quadratic programming (QP) in the unit simplex, namely, to find $\min_{p \in \Delta} p^\top b + p^\top A A^\top p$, where b is an n -vector. This relationship, and the generalization of our MEB algorithm to QP in the simplex, is discussed in §3.3; for more general background on QP in the simplex, and related problems, see for example [Cla08].

1.1 Related work

Perhaps the most closely related work is that of Grigoriadis and Khachiyan [GK95], who showed how to approximately solve a zero-sum game up to additive precision ε in time $\tilde{O}(\varepsilon^{-2}(n + d))$, where the game matrix is $n \times d$. This problem is analogous to ours, and our algorithm is similar in structure to theirs, but where we minimize over $p \in \Delta$ and maximize over $x \in \mathbb{B}$, their optimization has not only p but also x in a unit simplex.

Their algorithm (and ours) relies on sampling based on x and p , to estimate inner products $x^\top v$ or $p^\top w$ for vectors v and w that are rows or columns of A . For a vector $p \in \Delta$, this estimation is easily done by returning w_i with probability p_i .

For vectors $x \in \mathbb{B}$, however, the natural estimation technique is to pick i with probability x_i^2 , and return v_i/x_i . The estimator from this ℓ_2 sample is less well-behaved, since it is unbounded, and can have a high variance. While ℓ_2 sampling has been used in streaming applications [MW10], it has not previously found applications in optimization due to this high variance problem.

Indeed, it might seem surprising that sublinearity is at all possible, given that the correct classifier might be determined by very few examples, as shown in figure 2. It thus seems necessary to go over all examples at least once, instead of looking at noisy estimates based on sampling.

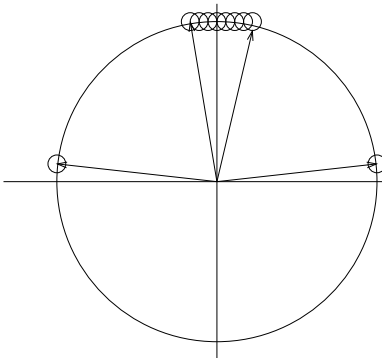


Figure 2: The optimum x_* is determined by the vectors near the horizontal axis.

However, as we show, in our setting there is a version of the fundamental Multiplicative Weights (MW) technique that can cope with unbounded updates, and for which the variance of ℓ_2 -sampling is manageable. In our version of MW, the multiplier associated with a value z is quadratic in z , in contrast to the more standard multiplier that is exponential in z ; while the latter is a fundamental building block in approximate optimization algorithms, as discussed by Plotkin *et al.* [PST91], in our setting such exponential updates can lead to a very expensive $d^{\Omega(1)}$ iterations.

We analyze MW from the perspective of on-line optimization, and show that our version of MW has low expected regret given only that the random updates have the variance bounds provable for ℓ_2 sampling. We also use another technique from on-line optimization, a gradient descent variant which is better suited for the ball.

For the special case of zero-sum games in which the entries are all non-negative (this is equivalent to packing and covering linear programs), Koufogiannakis and Young [KY07] give a sublinear-time algorithm which returns a *relative* approximation in time $\tilde{O}(\varepsilon^{-2}(n+d))$. Our lower bounds show that a similar relative approximation bound for sublinear algorithms is impossible for general classification, and hence general linear programming.

2 Linear Classification and the Perceptron

Before our algorithm, some reminders and further notation: $\Delta \subset \mathbb{R}^n$ is the unit simplex $\{p \in \mathbb{R}^n \mid p_i \geq 0, \sum_i p_i = 1\}$, $\mathbb{B} \subset \mathbb{R}^d$ is the Euclidean unit ball, and the unsubscripted $\|x\|$ denotes the Euclidean norm $\|x\|_2$. The n -vector, all of whose entries are one, is denoted by $\mathbf{1}_n$.

The i 'th row of the input matrix A is denoted A_i , although a vector is a column vector unless otherwise indicated. The i 'th coordinate of vector v is denoted $v(i)$. For a vector v , we let v^2 denote the vector whose coordinates have $v^2(i) \equiv v(i)^2$ for all i .

2.1 The Sublinear Perceptron

Our sublinear perceptron algorithm is given in Figure 1. The algorithm maintains a vector $w_t \in \mathbb{R}^n$, with nonnegative coordinates, and also $p_t \in \Delta$, which is w_t scaled to have unit ℓ_1 norm. A vector $y_t \in \mathbb{R}^d$ is maintained also, and x_t which is y_t scaled to have Euclidean norm no larger than one. These normalizations are done on line 4.

In lines 5 and 6, the algorithm is updating y_t by adding a row of A randomly chosen using p_t . This is a randomized version of *Online Gradient Descent* (OGD); due to the random choice of i_t , A_{i_t} is an unbiased estimator of $p_t^\top A$, which is the gradient of $p_t^\top A y$ with respect to y .

In lines 7 through 12, the algorithm is updating w_t using a column j_t of A randomly chosen based on x_t , and also using the value $x_t(j_t)$. This is a version of the Multiplicative Weights (MW) technique for online optimization in the unit simplex, where v_t is an unbiased estimator of $A x_t$, the gradient of $p^\top A x_t$ with respect to p .

Actually, v_t is not unbiased, after the clip operation: for $z, V \in \mathbb{R}$, $\text{clip}(z, V) \equiv \min\{V, \max\{-V, z\}\}$, and our analysis is helped by clipping the entries of v_t ; we show that the resulting slight bias is not harmful.

As discussed in §1.1, the sampling used to choose j_t (and update p_t) is ℓ_2 -sampling, and that for i_t , ℓ_1 -sampling. These techniques, which can be regarded as special cases of an ℓ_p -sampling technique, for $p \in [1, \infty)$, yield unbiased estimators of vector dot products. It is important for us also that ℓ_2 -sampling has a variance bound here; in particular, for each relevant i and t ,

$$\mathbf{E}[v_t(i)^2] \leq \|A_i\|^2 \|x_t\|^2 \leq 1. \tag{2}$$

First we note the running time.

Theorem 2.1. *The sublinear perceptron takes $O(\varepsilon^{-2} \log n)$ iterations, with a total running time of $O(\varepsilon^{-2}(n+d) \log n)$.*

Proof. The algorithm iterates $T = O(\frac{\log n}{\varepsilon^2})$ times. Each iteration requires:

1. One ℓ_2 sample per iterate, which takes $O(d)$ time using known data structures.

Algorithm 1 Sublinear Perceptron

1: Input: $\varepsilon > 0$, $A \in \mathbb{R}^{n \times d}$ with $A_i \in \mathbb{B}$ for $i \in [n]$.
 2: Let $T \leftarrow 200^2 \varepsilon^{-2} \log n$, $y_1 \leftarrow 0$, $w_1 \leftarrow \mathbf{1}_n$,
 $\eta \leftarrow \frac{1}{100} \sqrt{\frac{\log n}{T}}$.
 3: **for** $t = 1$ to T **do**
 4: $p_t \leftarrow \frac{w_t}{\|w_t\|_1}$, $x_t \leftarrow \frac{y_t}{\max\{1, \|y_t\|\}}$.
 5: Choose $i_t \in [n]$ by $i_t \leftarrow i$ with prob. $p_t(i)$.
 6: $y_{t+1} \leftarrow y_t + \frac{1}{\sqrt{2T}} A_{i_t}$
 7: Choose $j_t \in [d]$ by
 $j_t \leftarrow j$ with probability $x_t(j)^2 / \|x_t\|^2$.
 8: **for** $i \in [n]$ **do**
 9: $\tilde{v}_t(i) \leftarrow A_i(j_t) \|x_t\|^2 / x_t(j_t)$
 10: $v_t(i) \leftarrow \text{clip}(\tilde{v}_t(i), 1/\eta)$
 11: $w_{t+1}(i) \leftarrow w_t(i)(1 - \eta v_t(i) + \eta^2 v_t(i)^2)$
 12: **end for**
 13: **end for**
 14: **return** $\bar{x} = \frac{1}{T} \sum_t x_t$

2. Sampling $i_t \in_R p_t$ which takes $O(n)$ time.

3. The update of x_t and p_t , which takes $O(n + d)$ time.

The total running time is $O(\varepsilon^{-2}(n + d) \log n)$. □

Next we analyze the output quality. The proof uses new tools from regret minimization and sampling that are the building blocks of most of our upper bound results.

Let us first state the MW algorithm used in all our algorithms.

Definition 2.2 (MW algorithm). Consider a sequence of vectors $q_1, \dots, q_T \in \mathbb{R}^n$. The *Multiplicative Weights* (MW) algorithm is as follows. Let $w_1 \leftarrow \mathbf{1}_n$, and for $t \geq 1$,

$$p_t \leftarrow w_t / \|w_t\|_1, \tag{3}$$

and for $0 < \eta \in \mathbb{R}$

$$w_{t+1}(i) \leftarrow w_t(i)(1 - \eta q_t(i) + \eta^2 q_t(i)^2), \tag{4}$$

The following is a key lemma, which proves a novel bound on the regret of the MW algorithm above, suitable for the case where the losses are random variables with bounded variance. This is proven below, after a concentration lemma, and the main theorem and its proof.

Lemma 2.3 (Variance MW Lemma). *The MW algorithm satisfies*

$$\begin{aligned}
 \sum_{t \in [T]} p_t^\top q_t &\leq \min_{i \in [n]} \sum_{t \in [T]} \max\{q_t(i), -\frac{1}{\eta}\} \\
 &\quad + \frac{\log n}{\eta} + \eta \sum_{t \in [T]} p_t^\top q_t^2.
 \end{aligned}$$

The following three lemmas give concentration bounds on our random variables from their expectations. The first two are based on standard martingale analysis, and the last is a simple Markov application. The proofs are deferred to Appendix B.

Lemma 2.4. For $\eta \leq \sqrt{\frac{\log n}{10T}}$, with probability at least $1 - O(1/n)$,

$$\max_i \sum_{t \in [T]} [v_t(i) - A_i x_t] \leq 90\eta T.$$

Lemma 2.5. For $\eta \leq \sqrt{\frac{\log n}{10T}}$, with probability at least $1 - O(1/n)$, it holds that $\left| \sum_{t \in [T]} A_i x_t - \sum_t p_t^\top v_t \right| \leq 100\eta T$.

Lemma 2.6. With probability at least $1 - \frac{1}{4}$, it holds that $\sum_t p_t^\top v_t^2 \leq 8T$.

Theorem 2.7 (Main Theorem). With probability $1/2$, the sublinear perceptron returns a solution \bar{x} that is an ε -approximation.

Proof. First we use the regret bounds for lazy gradient descent to lower bound $\sum_{t \in [T]} A_i x_t$, next we get an upper bound for that quantity using the Weak Regret lemma above, and then we combine the two.

By definition, $A_i x^* \geq \sigma$ for all $i \in [n]$, and so, using the bound of Lemma A.2,

$$T\sigma \leq \max_{x \in \mathbb{B}} \sum_{t \in [T]} A_i x \leq \sum_{t \in [T]} A_i x_t + 2\sqrt{2T}, \quad (5)$$

or rearranging,

$$\sum_{t \in [T]} A_i x_t \geq T\sigma - 2\sqrt{2T}. \quad (6)$$

Now we turn to the MW part of our algorithm. By the Weak Regret Lemma 2.3, and using the clipping of $v_t(i)$,

$$\sum_{t \in [T]} p_t^\top v_t \leq \min_{i \in [n]} \sum_{t \in [T]} v_t(i) + (\log n)/\eta + \eta \sum_{t \in [T]} p_t^\top v_t^2.$$

By Lemma 2.4 above, with high probability, for any $i \in [n]$,

$$\sum_{t \in [T]} A_i x_t \geq \sum_{t \in [T]} v_t(i) - 90\eta T,$$

so that with high probability

$$\begin{aligned} \sum_{t \in [T]} p_t^\top v_t &\leq \min_{i \in [n]} \sum_{t \in [T]} A_i x_t + (\log n)/\eta \\ &\quad + \eta \sum_{t \in [T]} p_t^\top v_t^2 + 90T\eta. \end{aligned} \quad (7)$$

Combining (6) and (7) we get

$$\begin{aligned} \min_{i \in [n]} \sum_{t \in [T]} A_i x_t &\geq -(\log n)/\eta - \eta \sum_{t \in [T]} p_t^\top v_t^2 - 90T\eta \\ &\quad + T\sigma - 2\sqrt{2T} - \left| \sum_{t \in [T]} p_t^\top v_t - \sum_{t \in [T]} A_i x_t \right| \end{aligned}$$

By Lemmas 2.5, 2.6 we have w.p at least $\frac{3}{4} - O(\frac{1}{n}) \geq \frac{1}{2}$

$$\begin{aligned} \min_{i \in [n]} \sum_{t \in [T]} A_i x_t &\geq -(\log n)/\eta - 8\eta T - 90T\eta + T\sigma - 2\sqrt{2T} - 100\eta T \\ &\geq T\sigma - \frac{\log n}{\eta} - 200\eta T. \end{aligned}$$

Dividing through by T , and using our choice of η , we have $\min_i A_i \bar{x} \geq \sigma - \varepsilon/2$ w.p. at least $1/2$ as claimed. \square

Proof of Lemma 2.3, Weak Regret. We first show an upper bound on $\log \|w_{T+1}\|_1$, then a lower bound, and then relate the two.

From (4) and (3) we have

$$\begin{aligned} \|w_{t+1}\|_1 &= \sum_{i \in [n]} w_{t+1}(i) \\ &= \sum_{i \in [n]} p_t(i) \|w_t\|_1 (1 - \eta q_t(i) + \eta^2 q_t(i)^2) \\ &= \|w_t\|_1 (1 - \eta p_t^\top q_t + \eta^2 p_t^\top q_t^2). \end{aligned}$$

This implies by induction on t , and using $1 + z \leq \exp(z)$ for $z \in \mathbb{R}$, that

$$\log \|w_{T+1}\|_1 = \log n + \sum_{t \in [T]} \log(1 - \eta p_t^\top q_t + \eta^2 p_t^\top q_t^2) \leq \log n - \sum_{t \in [T]} \eta p_t^\top q_t + \eta^2 p_t^\top q_t^2. \quad (8)$$

Now for the lower bound. From (4) we have by induction on t that

$$w_{T+1}(i) = \prod_{t \in [T]} (1 - \eta q_t(i) + \eta^2 q_t(i)^2),$$

and so

$$\begin{aligned} \log \|w_{T+1}\|_1 &= \log \left[\sum_{i \in [n]} \prod_{t \in [T]} (1 - \eta q_t(i) + \eta^2 q_t(i)^2) \right] \\ &\geq \log \left[\max_{i \in [n]} \prod_{t \in [T]} (1 - \eta q_t(i) + \eta^2 q_t(i)^2) \right] \\ &= \max_{i \in [n]} \sum_{t \in [T]} \log(1 - \eta q_t(i) + \eta^2 q_t(i)^2) \\ &\geq \max_{i \in [n]} \sum_{t \in [T]} [\min\{-\eta q_t(i), 1\}], \end{aligned}$$

where the last inequality uses the fact that $1 + z + z^2 \geq \exp(\min\{z, 1\})$ for all $z \in \mathbb{R}$.

Putting this together with the upper bound (8), we have

$$\max_{i \in [n]} \sum_{t \in [T]} [\min\{-\eta q_t(i), 1\}] \leq \log n - \sum_{t \in [T]} \eta p_t^\top q_t + \eta^2 p_t^\top q_t^2,$$

Changing sides

$$\begin{aligned} \sum_{t \in [T]} \eta p_t^\top q_t &\leq -\max_{i \in [n]} \sum_{t \in [T]} [\min\{-\eta q_t(i), 1\}] + \log n + \eta^2 p_t^\top q_t^2, \\ &= \min_{i \in [n]} \sum_{t \in [T]} [\max\{\eta q_t(i), -1\}] + \log n + \eta^2 p_t^\top q_t^2, \end{aligned}$$

and the lemma follows, dividing through by η . \square

Corollary 2.8 (Dual solution). *The vector $\bar{p} \equiv \sum_t e_{i_t}/T$ is, with probability $1/2$, an $O(\varepsilon)$ -approximate dual solution.*

Proof. Observing in (5) that the middle expression $\max_{x \in \mathbb{B}} \sum_{t \in [T]} A_{i_t} x$ is equal to $T \max_{x \in \mathbb{B}} \bar{p}^\top A x$, we have $T \max_{x \in \mathbb{B}} \bar{p}^\top A x \leq \sum_{t \in [T]} A_{i_t} x_t + 2\sqrt{2T}$, or changing sides,

$$\sum_{t \in [T]} A_{i_t} x_t \geq T \max_{x \in \mathbb{B}} \bar{p}^\top A x - 2\sqrt{2T}$$

Recall from (7) that with high probability,

$$\sum_{t \in [T]} p_t^\top v_t \leq \min_{i \in [n]} \sum_{t \in [T]} A_{i_t} x_t + (\log n)/\eta + \eta \sum_{t \in [T]} p_t^\top v_t^2 + 90T\eta. \quad (9)$$

Following the proof of the main Theorem, we combine both inequalities and use Lemmas 2.5,2.6, such that with probability at least $\frac{1}{2}$:

$$\begin{aligned} T \max_{x \in \mathbb{B}} \bar{p}^\top A x &\leq \min_{i \in [n]} \sum_{t \in [T]} A_{i_t} x_t + (\log n)/\eta + \eta \sum_{t \in [T]} p_t^\top v_t^2 + 90T\eta + 2\sqrt{2T} + \left| \sum_{t \in [T]} p_t^\top v_t - \sum_{t \in [T]} A_{i_t} x_t \right| \\ &\leq T\sigma + O(\sqrt{T \log n}) \end{aligned}$$

Dividing through by T we have with probability at least $\frac{1}{2}$ that $\max_{x \in \mathbb{B}} \bar{p}^\top A x \leq \sigma + O(\varepsilon)$ for our choice of T and η . \square

2.2 High Success Probability and Las Vegas

Given two vectors $u, v \in \mathbb{B}$, we have seen that a single ℓ_2 -sample is an unbiased estimator of their inner product with variance at most one. Averaging $\frac{1}{\varepsilon^2}$ such samples reduces the variance to ε^2 , which reduces the standard deviation to ε . Repeating $O(\log \frac{1}{\delta})$ such estimates, and taking the median, gives an estimator denoted $X_{\varepsilon, \delta}$, which satisfies, via a Chernoff bound:

$$\Pr[|X_{\varepsilon, \delta} - v^\top u| > \varepsilon] \leq \delta$$

As an immediate corollary of this fact we obtain:

Corollary 2.9. *There exists a randomized algorithm that with probability $1 - \delta$, successfully determines whether a given hyperplane with normal vector $x \in \mathbb{B}$, together with an instance of linear classification and parameter $\sigma > 0$, is an ε -approximate solution. The algorithm runs in time $O(d + \frac{n}{\varepsilon^2} \log \frac{n}{\delta})$.*

Proof. Let $\delta' = \delta/n$. Generate the random variable $X_{\varepsilon, \delta'}$ for each inner product pair $\langle x, A_i \rangle$, and return true if and only if $X_{\varepsilon, \delta'} \geq \sigma - \varepsilon$ for each pair. By the observation above and taking union bound over all n inner products, with probability $1 - \delta$ the estimate $X_{\varepsilon, \delta'}$ was ε -accurate for all inner-product pairs, and hence the algorithm returned a correct answer.

The running time includes preprocessing of x in $O(d)$ time, and n inner-product estimates, for a total of $O(d + \frac{n}{\varepsilon^2} \log \frac{n}{\delta})$. \square

Hence, we can amplify the success probability of Algorithm 1 to $1 - \delta$ for any $\delta > 0$ albeit incurring additional poly-log factors in running time:

Corollary 2.10 (High probability). *There exists a randomized algorithm that with probability $1 - \delta$ returns an ε -approximate solution to the linear classification problem, and runs in expected time $O(\frac{n+d}{\varepsilon^2} \log \frac{n}{\delta})$.*

Proof. Run Algorithm 1 for $\log_2 \frac{1}{\delta}$ times to generate that many candidate solutions. By Theorem 2.7, at least one candidate solution is an ε -approximate solution with probability at least $1 - 2^{-\log_2 \frac{1}{\delta}} = 1 - \delta$.

For each candidate solution apply the verification procedure above with success probability $1 - \delta^2 \geq 1 - \frac{\delta}{\log \frac{1}{\delta}}$, and all verifications will be correct again with probability at least $1 - \delta$. Hence, both events hold with probability at least $1 - 2\delta$. The result follows after adjusting constants.

The worst-case running time comes to $O(\frac{n+d}{\varepsilon^2} \log \frac{n}{\delta} \log \frac{1}{\delta})$. However, we can generate the candidate solutions and verify them one at a time, rather than all at once. The expected number of candidates we need to generate is constant. \square

It is also possible to obtain an algorithm that never errs:

Corollary 2.11 (Las Vegas Version). *After $O(\varepsilon^{-2} \log n)$ iterations, the sublinear perceptron returns a solution that with probability $1/2$ can be verified in $O(M)$ time to be ε -approximate. Thus with expected $O(1)$ repetitions, and a total of expected $O(M + \varepsilon^{-2}(n + d) \log n)$ work, a verified ε -approximate solution can be found.*

Proof. We have

$$\min_i A_i \bar{x} \leq \sigma \leq \|\bar{p}^\top A\|,$$

and so if

$$\min_i A_i \bar{x} \geq \|\bar{p}^\top A\| - \varepsilon, \tag{10}$$

then \bar{x} is an ε -approximate solution, and \bar{x} will pass this test if it and \bar{p} are $(\varepsilon/2)$ -approximate solutions, and the same for \bar{p} .

Thus, running the algorithm for a constant factor more iterations, so that with probability $1/2$, \bar{x} and \bar{p} are both $(\varepsilon/2)$ -approximate solutions, it can be verified that both are ε -approximate solutions. \square

2.3 Further Optimizations

The regret of OGD as given in Lemma A.2 is smaller than the dual strategy of random MW. We can take advantage of this and improve the running time slightly, by replacing line [6] of the sublinear algorithm with the line shown below.

This has the effect of increasing the regret of the primal online algorithm by a log n factor, which does not hurt the number of iterations required to converge, since the overall regret is dominated by that of the MW algorithm.

[6'] With probability $\frac{1}{\log T}$, let $y_{t+1} \leftarrow y_t + \frac{1}{2\sqrt{T}}A_{i_t}$ (else do nothing).

Since the primal solution x_t is not updated in every iteration, we improve the running time slightly to

$$O(\varepsilon^{-2} \log n(n + d/(\log 1/\varepsilon + \log \log n))).$$

We use this technique to greater effect for the MEB problem, where it is discussed in more detail.

2.4 Implications in the PAC model

Consider the “separable” case of hyperplane learning, in which there exists a hyperplane classifying all data points correctly. It is well known that the concept class of hyperplanes in d dimensions with margin σ has effective dimension at most $\min\{d, \frac{1}{\sigma^2}\} + 1$. Consider the case in which the margin is significant, i.e. $\frac{1}{\sigma^2} < d$. PAC learning theory implies that the number of examples needed to attain generalization error of δ is $O(\frac{1}{\sigma^2\delta})$.

Using the method of online to batch conversion (see [CBCG04]), and applying the online gradient decent algorithm, it is possible to obtain δ generalization error in time $O(\frac{d}{\sigma^2\delta})$ time, by going over the data once and performing a gradient step on each example.

Our algorithm improves upon this running time bound as follows: we use the sublinear perceptron to compute a $\sigma/2$ -approximation to the best hyperplane over the test data, where the number of examples is taken to be $n = O(\frac{1}{\sigma^2\delta})$ (in order to obtain δ generalization error). As shown previously, the total running time amounts to $\tilde{O}(\frac{\frac{1}{\sigma^2\delta} + d}{\sigma^4\delta}) = O(\frac{1}{\sigma^4\delta} + \frac{d}{\sigma^2})$.

This improves upon standard methods by a factor of $\tilde{O}(\sigma^2d)$, which is always an improvement by our initial assumption on σ and d .

3 Strongly convex problems: MEB and SVM

3.1 Minimum Enclosing Ball

In the Minimum Enclosing Ball problem the input consists of a matrix $A \in \mathbb{R}^{n \times d}$. The rows are interpreted as vectors and the problem is to find a vector $x \in \mathbb{R}^d$ such that

$$x_* \equiv \operatorname{argmin}_{x \in \mathbb{R}^d} \max_{i \in [n]} \|x - A_i\|^2$$

We further assume for this problem that all vectors A_i have Euclidean norm at most one. Denote by $\sigma = \max_{i \in [n]} \|x - A_i\|^2$ the radius of the optimal ball, and we say that a solution is ε -approximate if the ball it generates has radius at most $\sigma + \varepsilon$.

As in the case of linear classification, to obtain tight running time bounds we use a primal-dual approach; the algorithm is given below.

(This is a “conceptual” version of the algorithm: in the analysis of the running time, we use the fact that we can batch together the updates for w_t over the iterations for which x_t does not change.)

Theorem 3.1. *Algorithm 2 runs in $O(\frac{\log n}{\varepsilon^2})$ iterations, with a total expected running time of*

$$\tilde{O}\left(\frac{n}{\varepsilon^2} + \frac{d}{\varepsilon}\right),$$

and with probability $1/2$, returns an ε -approximate solution.

Algorithm 2 Sublinear Primal-Dual MEB

- 1: Input: $\varepsilon > 0$, $A \in \mathbb{R}^{n \times d}$ with $A_i \in \mathbb{B}$ for $i \in [n]$ and $\|A_i\|$ known.
 - 2: Let $T \leftarrow \Theta(\varepsilon^{-2} \log n)$, $y_1 \leftarrow \mathbf{0}$, $w_1 \leftarrow \mathbf{1}$, $\eta \leftarrow \sqrt{(\log n)/T}$, $\alpha \leftarrow \frac{\log T}{\sqrt{T \log n}}$.
 - 3: **for** $t = 1$ to T **do**
 - 4: $p_t \leftarrow \frac{w_t}{\|w_t\|_1}$
 - 5: Choose $i_t \in [n]$ by $i_t \leftarrow i$ with probability $p_t(i)$.
 - 6: With probability α , update $y_{t+1} \leftarrow y_t + A_{i_t}$, $x_{t+1} \leftarrow \frac{y_{t+1}}{t}$. (else do nothing)
 - 7: Choose $j_t \in [d]$ by $j_t \leftarrow j$ with probability $x_t(j)^2 / \|x_t\|^2$.
 - 8: **for** $i \in [n]$ **do**
 - 9: $\tilde{v}_t(i) \leftarrow -2A_i(j_t)\|x_t\|^2/x_t(j_t) + \|A_i\|^2 + \|x_t\|^2$.
 - 10: $v_t(i) \leftarrow \text{clip}(\tilde{v}_t(i), \frac{1}{\eta})$.
 - 11: $w_{t+1}(i) \leftarrow w_t(i)(1 + \eta v_t(i) + \eta^2 v_t(i)^2)$.
 - 12: **end for**
 - 13: **end for**
 - 14: **return** $\bar{x} = \frac{1}{T} \sum_t x_t$
-

Proof. Except for the running time analysis, the proof of this theorem is very similar to that of Theorem 2.7, where we take advantage of a tighter regret bound for strictly convex loss functions in the case of MEB, for which the OGD algorithm with a learning rate of $\frac{1}{t}$ is known to obtain a tighter regret bound of $O(\log T)$ instead of $O(\sqrt{T})$. For presentation, we use asymptotic notation rather than computing the exact constants (as done for the linear classification problem).

Let $f_t(x) = \|x - A_{i_t}\|^2$. Notice that $\arg \min_{x \in \mathbb{B}} \sum_{\tau=1}^t f_\tau(x) = \frac{\sum_{\tau=1}^t A_{i_\tau}}{t}$. By Lemma A.5 such that $f_t(x) = \|x - A_{i_t}\|^2$, with $G \leq 2$ and $H = 2$, and x^* being the solution to the instance, we have

$$\mathbf{E}_{\{c_t\}} \left[\sum_t \|x_t - A_{i_t}\|^2 \right] \leq \mathbf{E}_{\{c_t\}} \left[\sum_t \|x^* - A_{i_t}\|^2 \right] + \frac{4}{\alpha} \log T \leq T\sigma + \frac{4}{\alpha} \log T, \quad (11)$$

where σ is the squared MEB radius. Here the expectation is taken only over the random coin tosses for updating x_t , denoted c_t , and holds for any outcome of the indices i_t sampled from p_t and the coordinates j_t used for the ℓ_2 sampling.

Now we turn to the MW part of our algorithm. By the Weak Regret Lemma 2.3, using the clipping of $v_t(i)$, and reversing inequalities to account for the change of sign, we have

$$\sum_{t \in [T]} p_t^\top v_t \geq \max_{i \in [n]} \sum_{t \in [T]} v_t(i) - O\left(\frac{\log n}{\eta} + \eta \sum_{t \in [T]} p_t^\top v_t^2\right).$$

Using Lemmas B.4, B.5 with high probability

$$\forall i \in [n]. \sum_{t \in [T]} v_t(i) \geq \sum_{t \in [T]} \|A_i - x_t\|^2 - O(\eta T),$$

$$\left| \sum_{t \in [T]} \|x_t - A_{i_t}\|^2 - \sum_{t \in [T]} p_t^\top v_t \right| = O(\eta T).$$

Plugging these two facts in the previous inequality we have w.h.p

$$\sum_{t \in [T]} \|x_t - A_{i_t}\|^2 \geq \max_{i \in [n]} \sum_{t \in [T]} \|A_i - x_t\|^2 - O\left(\frac{\log n}{\eta} + \eta \sum_{t \in [T]} p_t^\top v_t^2 + T\eta\right).$$

This holds w.h.p over the random choices of $\{i_t, j_t\}$, and irrespective of the coin tosses $\{c_t\}$. Hence, we can take expectations w.r.t $\{c_t\}$, and obtain

$$\mathbf{E}_{\{c_t\}}\left[\sum_{t \in [T]} \|x_t - A_{i_t}\|^2\right] \geq \mathbf{E}_{\{c_t\}}\left[\max_{i \in [n]} \sum_{t \in [T]} \|A_i - x_t\|^2\right] - O\left(\frac{\log n}{\eta} + \eta \sum_{t \in [T]} p_t^\top v_t^2 + T\eta\right). \quad (12)$$

Combining with equation (11), we obtain that w.h.p. over the random variables $\{i_t, j_t\}$

$$T\sigma + \frac{4}{\alpha} \log T \geq \mathbf{E}_{\{c_t\}}\left[\max_{i \in [n]} \sum_{t \in [T]} \|x_t - A_i\|^2\right] - O\left(\frac{\log n}{\eta} + \eta \sum_{t \in [T]} p_t^\top v_t^2 + T\eta\right)$$

Rearranging and using Lemma B.8, we have w.p. at least $\frac{1}{2}$

$$\mathbf{E}_{\{c_t\}}\left[\max_{i \in [n]} \sum_{t \in [T]} \|x_t - A_i\|^2\right] \leq O\left(T\sigma + \frac{\log T}{\alpha} + \frac{\log n}{\eta} + T\eta\right)$$

Dividing through by T and applying Jensen's inequality, we have

$$\mathbf{E}\left[\max_j \|\bar{x} - A_j\|^2\right] \leq \frac{1}{T} \mathbf{E}\left[\max_{i \in [n]} \sum_{t \in [T]} \|x_t - A_i\|^2\right] \leq O\left(\sigma + \frac{\log T}{T\alpha} + \frac{\log n}{T\eta} + \eta\right).$$

Optimizing over the values of α , η , and T , this implies that the expected error is $O(\varepsilon)$, and so using Markov's inequality, \bar{x} is a $O(\varepsilon)$ -approximate solution with probability at least $1/2$.

Running time The algorithm above consists of $T = O\left(\frac{\log n}{\varepsilon^2}\right)$ iterations. Naively, this would result in the same running time as for linear classification. Yet notice that x_t changes only an expected αT times, and only then do we perform an $O(d)$ operation. The expected number of iterations in which x_t changes is $\alpha T \leq 16\varepsilon^{-1} \log T$, and so the running time is

$$O\left(\varepsilon^{-1}(\log T) \cdot d + \frac{\log n}{\varepsilon^2} \cdot n\right) = \tilde{O}(\varepsilon^{-2}n + \varepsilon^{-1}d).$$

□

The following Corollary is a direct analogue of Corollary 2.8.

Corollary 3.2 (Dual solution). *The vector $\bar{p} \equiv \sum_t e_{i_t}/T$ is, with probability $1/2$, an $O(\varepsilon)$ -approximate dual solution.*

3.2 High Success Probability and Las Vegas

As for linear classification, we can amplify the success probability of Algorithm 2 to $1 - \delta$ for any $\delta > 0$ albeit incurring additional poly-log factors in running time.

Corollary 3.3 (MEB high probability). *There exists a randomized algorithm that with probability $1 - \delta$ returns an ε -approximate solution to the MEB problem, and runs in expected time $\tilde{O}\left(\frac{n}{\varepsilon^2} \log \frac{n}{\varepsilon\delta} + \frac{d}{\varepsilon} \log \frac{1}{\varepsilon}\right)$. There is also a randomized algorithm that returns an ε -approximate solution in $\tilde{O}\left(M + \frac{n}{\varepsilon^2} + \frac{d}{\varepsilon}\right)$ time.*

Proof. We can estimate the distance between two points in \mathbb{B} in $O(\varepsilon^{-2} \log(1/\delta))$ time, with error at most ε and failure probability at most δ , using the dot product estimator described in §2.2. Therefore we can estimate the maximum distance of a given point to every input point in $O(n\varepsilon^{-2} \log(n/\delta))$ time, with error at most ε and failure probability at most δ . This distance is $\sigma - \varepsilon$, where σ is the optimal radius attainable, w.p. $1 - \delta$.

Because Algorithm 2 yields an ε -dual solution with probability $1/2$, we can use this solution to verify that the radius of any possible solution to the farthest point is at least $\sigma - \varepsilon$.

So, to obtain a solution as described in the lemma statement, run Algorithm 2, and verify that it yields an ε -approximation, using this approximate dual solution; with probability $1/2$, this gives a verified ε -approximation. Keep trying until this succeeds, in an expected 2 trials.

For a Las Vegas algorithm, we simply apply the same scheme, but verify the distances exactly. \square

3.3 Convex Quadratic Programming in the Simplex

We can extend our approach to problems of the form

$$\min_{p \in \Delta} p^\top b + p^\top A A^\top p, \quad (13)$$

where $b \in \mathbb{R}^n$, $A \in \mathbb{R}^{n \times d}$, and Δ is, as usual, the unit simplex in \mathbb{R}^n . As is well known, and as we partially review below, this problem includes the MEB problem, margin estimation as for hard margin support vector machines, the L_2 -SVM variant of support vector machines, the problem of finding the shortest vector in a polytope, and others.

Applying $\|v - x\|^2 = v^\top v + x^\top x - 2v^\top x \geq 0$ with $v \leftarrow A^\top p$, we have

$$\max_{x \in \mathbb{R}^d} 2p^\top A x - \|x\|^2 = p^\top A A^\top p, \quad (14)$$

with equality at $x = A^\top p$. Thus (13) can be written as

$$\min_{p \in \Delta} \max_{x \in \mathbb{R}^d} p^\top (b + 2A x - \mathbf{1}_n \|x\|^2). \quad (15)$$

The *Wolfe dual* of this problem exchanges the max and min:

$$\max_{x \in \mathbb{R}^d} \min_{p \in \Delta} p^\top (b + 2A x - \mathbf{1}_n \|x\|^2). \quad (16)$$

Since

$$\min_{p \in \Delta} p^\top (b + 2A x - \mathbf{1}_n \|x\|^2) = \min_i b(i) + 2A_i x - \|x\|^2, \quad (17)$$

with equality when $p_i = 0$ if \hat{i} is not a minimizer, the dual can also be expressed as

$$\max_{x \in \mathbb{R}^d} \min_i b(i) + 2A_i x - \|x\|^2 \quad (18)$$

By the two relations (14) and (17) used to derive the dual problem from the primal, we have immediately the *weak duality* condition that the objective function of the dual (18) is always no more than the objective function value of the primal (13). The strong duality condition, that the two problems take the same optimal value, also holds here; indeed, the optimum x_* also solves (14), and the optimal p_* also solves (17).

To generalize Algorithm 2, we make v_t an unbiased estimator of $b + 2Ax_t - \mathbf{1}_n \|x_t\|^2$, and set x_{t+1} to be the minimizer of

$$\sum_{t' \in [t]} b(i_{t'}) + 2A_{i_{t'}} x_{t'} - \|x_{t'}\|^2,$$

namely, as with MEB, $y_{t+1} \leftarrow \sum_{t' \in [t]} A_{i_{t'}}$, and $x_{t+1} \leftarrow y_{t+1}/t$. (We also make some sign changes to account for the max-min formulation here, versus the min-max formulation used for MEB above.) This allows the use of Lemma A.4 for essentially the same analysis as for MEB; the gradient bound G and Hessian bound H are both at most 2, again assuming that all $A_i \in \mathbb{B}$.

MEB When the $b(i) \leftarrow -\|A_i\|^2$, we have

$$-\max_{x \in \mathbb{R}^d} \min_i b(i) + 2A_i x - \|x\|^2 = \min_{x \in \mathbb{R}^d} \max_i \|A_i\|^2 - 2A_i x + \|x\|^2 = \min_{x \in \mathbb{R}^d} \max_i \|x - A_i\|^2,$$

the objective function for the MEB problem.

Margin Estimation When $b \leftarrow 0$ in the primal problem (13), that problem is one of finding the shortest vector in the polytope $\{A^\top p \mid p \in \Delta\}$. Considering this case of the dual problem (18), for any given $x \in \mathbb{R}^d$ with $\min_i A_i x \leq 0$, the value of $\beta \in \mathbb{R}$ such that βx maximizes $\min_i 2A_i \beta x - \|\beta x\|^2$ is $\beta = 0$. On the other hand if x is such that $\min_i A_i x > 0$, the maximizing value β is $\beta = A_i x / \|x\|^2$, so that the solution of (18) also maximizes $\min_i (A_i x)^2 / \|x\|^2$. The latter is the square of the margin σ , which as before is the minimum distance of the points A_i to the hyperplane that is normal to x and passes through the origin.

Adapting Algorithm 2 for margin estimation, and with the slight changes needed for its analysis, we have that there is an algorithm taking $\tilde{O}(n/\epsilon^2 + d/\epsilon)$ time that finds $\bar{x} \in \mathbb{R}^d$ such that, for all $i \in [n]$,

$$2A_i \bar{x} - \|\bar{x}\|^2 \geq \sigma^2 - \epsilon.$$

When $\sigma^2 \leq \epsilon$, we don't appear to gain any useful information. However, when $\sigma^2 > \epsilon$, we have $\min_{i \in [n]} A_i \bar{x} > 0$, and so, by appropriate scaling of \bar{x} , we have \hat{x} such that

$$\hat{\sigma}^2 = \min_{i \in [n]} (A_i \hat{x})^2 / \|\hat{x}\|^2 = \min_{i \in [n]} 2A_i \hat{x} - \|\hat{x}\|^2 \geq \sigma^2 - \epsilon,$$

and so $\hat{\sigma} \geq \sigma - \epsilon/\sigma$. That is, letting $\epsilon \equiv \epsilon' \sigma$, if $\epsilon' \leq \sigma$, there is an algorithm taking $\tilde{O}(n/(\epsilon \sigma)^2 + d/\epsilon' \sigma)$ time that finds a solution \hat{x} with $\hat{\sigma} \geq \sigma - \epsilon'$.

4 A Generic Sublinear Primal-Dual Algorithm

We note that our technique above can be applied more broadly to any constrained optimization problem for which low-regret algorithms exist and low-variance sampling can be applied efficiently; that is, consider the general problem with optimum σ :

$$\max_{x \in \mathcal{K}} \min_i c_i(x) = \sigma. \tag{19}$$

Suppose that for the set \mathcal{K} and cost functions $c_i(x)$, there exists an iterative low regret algorithm, denoted LRA , with regret $R(T) = o(T)$. Let $T_\epsilon(LRA)$ be the smallest T such that $\frac{R(T)}{T} \leq \epsilon$. We denote by $x_{t+1} \leftarrow LRA(x_t, c)$ an invocation of this algorithm, when at state $x_t \in \mathcal{K}$ and the cost function c is observed.

Algorithm 3 Generic Sublinear Primal-Dual Algorithm

```

1: Let  $T \leftarrow \max\{T_\varepsilon(LRA), \frac{\log n}{\varepsilon^2}\}$  ,
    $x_1 \leftarrow LRA(\text{initial})$ ,  $w_1 \leftarrow \mathbf{1}_n$ ,  $\eta \leftarrow \frac{1}{100} \sqrt{\frac{\log n}{T}}$ .
2: for  $t = 1$  to  $T$  do
3:   for  $i \in [n]$  do
4:     Let  $v_t(i) \leftarrow \mathbf{Sample}(x_t, c_i)$ 
5:      $v_t(i) \leftarrow \text{clip}(\tilde{v}_t(i), 1/\eta)$ 
6:      $w_{t+1}(i) \leftarrow w_t(i)(1 - \eta v_t(i) + \eta^2 v_t(i)^2)$ 
7:   end for
8:    $p_t \leftarrow \frac{w_t}{\|w_t\|_1}$ ,
9:   Choose  $i_t \in [n]$  by  $i_t \leftarrow i$  with probability  $p_t(i)$ .
10:   $x_t \leftarrow LRA(x_{t-1}, c_{i_t})$ 
11: end for
12: return  $\bar{x} = \frac{1}{T} \sum_t x_t$ 

```

Let $\mathbf{Sample}(x, c)$ be a procedure that returns an unbiased estimate of $c(x)$ with variance at most one, that runs in constant time. Further assume $|c_i(x)| \leq 1$ for all $x \in K$, $i \in [n]$.

Applying the techniques of section 2 we can obtain the following generic lemma.

Lemma 4.1. *The generic sublinear primal-dual algorithm returns a solution x that with probability at least $\frac{1}{2}$ is an ε -approximate solution in $\max\{T_\varepsilon(LRA), \frac{\log n}{\varepsilon^2}\}$ iterations.*

Proof. First we use the regret bounds for LRA to lower bound $\sum_{t \in [T]} c_{i_t}(x_t)$, next we get an upper bound for that quantity using the Weak Regret Lemma, and then we combine the two in expectation.

By definition, $c_i(x^*) \geq \sigma$ for all $i \in [n]$, and so, using the LRA regret guarantee,

$$T\sigma \leq \max_{x \in \mathbb{B}} \sum_{t \in [T]} c_{i_t}(x) \leq \sum_{t \in [T]} c_{i_t}(x_t) + R(T), \quad (20)$$

or rearranging,

$$\sum_{t \in [T]} c_{i_t}(x_t) \geq T\sigma - R(T). \quad (21)$$

Now we turn to the MW part of our algorithm. By the Weak Regret Lemma 2.3, and using the clipping of $v_t(i)$,

$$\sum_{t \in [T]} p_t^\top v_t \leq \min_{i \in [n]} \sum_{t \in [T]} v_t(i) + (\log n)/\eta + \eta \sum_{t \in [T]} p_t^\top v_t^2.$$

Using Lemma B.4 and Lemma B.5, since the procedure \mathbf{Sample} is unbiased and has variance at most one, with high probability:

$$\forall i \in [n] \ , \ \sum_{t \in [T]} v_t(i) \leq \sum_{t \in [T]} c_i(x_t) + O(\eta T),$$

$$\left| \sum_{t \in [T]} c_{i_t}(x_t) - \sum_t p_t^\top v_t \right| = O(\eta T).$$

Plugging these two facts in the previous inequality we have w.h.p,

$$\sum_{t \in [T]} c_{i_t}(x_t) \leq \min_{i \in [n]} \sum_{t \in [T]} c_i(x_t) + O\left(\frac{\log n}{\eta} + \eta \sum_{t \in [T]} p_t^\top v_t^2 + \eta T\right) \quad (22)$$

Combining (21) and (22) we get w.h.p

$$\min_{i \in [n]} \sum_{t \in [T]} c_i(x_t) \geq -O\left(\frac{\log n}{\eta} + \eta T + \eta \sum_{t \in [T]} p_t^\top v_t^2\right) - R(T)$$

And via Lemma B.8 we have w.p. at least $\frac{1}{2}$ that

$$\min_{i \in [n]} \sum_{t \in [T]} c_i(x_t) \geq -O\left(\frac{\log n}{\eta} + \eta T\right) - R(T)$$

Dividing through by T , and using our choice of η , we have $\min_i c_i \bar{x} \geq \sigma - \varepsilon/2$ w.p. at least $1/2$ as claimed. \square

High-probability results can be obtained using the same technique as for linear classification.

4.1 More applications

The generic algorithm above can be used to derive the result of Grigoriadis and Khachiyan [GK95] on sublinear approximation of zero sum games with payoffs/losses bounded by one (up to polylogarithmic factors in running time). A zero sum game can be cast as the following min-max optimization problem:

$$\min_{x \in \Delta_d} \max_{i \in \Delta_n} A_i x$$

That is, the constraints are inner products with the rows of the game matrix. This is exactly the same as the linear classification problem, but the vectors x are taken from the convex set \mathcal{K} which is the simplex - or the set of all mixed strategies of the column player.

A low regret algorithm for the simplex is the multiplicative weights algorithm, which attains regret $R(T) \leq 2\sqrt{T \log n}$. The procedure **Sample**(x, A_i) to estimate the inner product $A_i x$ is much simpler than the one used for linear classification: we sample from the distribution x and return $A_i(j)$ w.p. $x(j)$. This has correct expectation and variance bounded by one (in fact, the random variable is always bounded by one). Lemma 4.1 then implies:

Corollary 4.2. *The sublinear primal-dual algorithm applied to zero sum games returns a solution x that with probability at least $\frac{1}{2}$ is an ε -approximate solution in $O\left(\frac{\log n}{\varepsilon^2}\right)$ iterations and total time $\tilde{O}\left(\frac{n+d}{\varepsilon^2}\right)$.*

Essentially any constrained optimization problem which has convex or linear constraints, and is over a simple convex body such as the ball or simplex, can be approximated in sublinear time using our method. The particular application to soft margin SVM, together with its practical significance, is explored in ongoing work with Nati Srebro.

5 A Semi-Streaming Implementation

In order to achieve space that is sublinear in d , we cannot afford to output a solution vector. We instead output both the cost of the solution, and a set of indices i_1, \dots, i_t for which the solution is a linear combination (that we know) of A_{i_1}, \dots, A_{i_t} . We note that all previous algorithms for these problems, even to achieve this notion of output, required $\Omega(d)$ space and/or $\Omega(nd)$ time, see, e.g., the references in [AS10].

We discuss the modifications to the sublinear primal-dual algorithm that need to be done for classification and minimum enclosing ball problems.

Our algorithm assumes it sees entire points at a time, i.e., it sees the entries of A row at a time, though the rows may be ordered arbitrarily. It relies on two streaming results about a d -dimensional vector x undergoing updates to its coordinates. We assume that each update is of the form (i, z) , where $i \in [d]$ is a coordinate of x and $z \in \{-P, -P + 1, \dots, P\}$ indicates that $x_i \leftarrow x_i + z$. The first is an efficient ℓ_2 -sketching algorithm of Thorup and Zhang. This algorithm allows for $(1 + \varepsilon)$ -approximation of $\|x\|_2$ with high probability using 1-pass, $\tilde{O}(\varepsilon^{-2})$ space, and time proportional to the length of the stream.

Theorem 5.1. ([TZ04]) *There is a 1-pass algorithm which outputs a $(1 \pm \varepsilon)$ -approximation to $\|x\|_2$ with probability $\geq 1 - \delta$ using $O(\varepsilon^{-2} \log(PdQ) \log 1/\delta)$ bits of space and $O(Q \log 1/\delta)$ time, where Q is the total number of updates in the stream.*

The second component is due to Monemizadeh and Woodruff [MW10]. We are given a stream of updates to a d -dimensional vector x , and want to output a random coordinate $I \in [d]$ for which for any $j \in [d]$, $\Pr[I = j] = \frac{|x_j|^2}{\|x\|_2^2}$. We also want the algorithm to return the value x_I . Such an algorithm is called an exact augmented ℓ_2 -Sampler. As shown in [MW10], an augmented ℓ_2 -Sampler with $O(\log d)$ space, $\tilde{O}(1)$ passes, and running time $\tilde{O}(Q)$ exists, where Q is the number of updates in the stream. This is what we use to ℓ_2 -sample from an iterate vector that we can only afford to represent implicitly.

Theorem 5.2. (Theorem 1.3 of [MW10]) *There is an $O(\log d)$ -pass exact augmented ℓ_2 -Sampler that uses $O(\log^5(Pd))$ bits of space and has running time $Q \log^{O(1)}(PdQ)$, where Q is the total number of updates in the stream. The algorithm fails with probability $\leq d^{-c}$ for an arbitrarily large constant $c > 0$.*

We maintain the indices i_t and j_t used in all $\tilde{O}(\varepsilon^{-2})$ iterations of the primal dual algorithm. Notice that in a single iteration t the same ℓ_2 -sample index j_t can be used for all n rows. While we cannot afford to remember the probabilities in the dual vector, we can store the values $\frac{\alpha_t}{x_t(j)}$, where α_t is a $(1 \pm \varepsilon)$ -approximation of $\|x_t\|_2^2$ which can be obtained using the Thorup-Zhang sketch. We also need such an approximation to $\|x_t\|$ to appropriately weight the rows used to do ℓ_2 -sampling (see below). Since we see rows (i.e., points) of A at a time, we can reconstruct the probability of each row in the dual vector on the fly in low space, and can use reservoir sampling to make the next choice of i_t . Then we use an augmented ℓ_2 -sampler to make the next choice of j_t , where we must ℓ_2 sample from a weighted sum of rows indexed by i_1, \dots, i_t in low space. We use the fact argued in §C We can show that the algorithm remains correct given the per-iteration rounding of the updates $v_t(i)$ to relative error μ , where μ is on the order of $\eta\varepsilon/T$. Throughout we round matrix entries to the nearest integer multiple of $\text{poly}(1/d)$ for a sufficiently large polynomial.

We implicitly represent the primal and dual vectors. At iteration t of the sublinear primal-dual algorithm, we have indices i_1, \dots, i_{t-1} of the sampled rows and indices j_1, \dots, j_t of the sampled columns for ℓ_2 -sampling (in a given iteration t , we use the same column j_t for ℓ_2 -sampling from all

rows). We maintain $\mu/2$ -approximations $\frac{1}{\tilde{x}_1(j_1)}, \dots, \frac{1}{\tilde{x}_{t-1}(j_{t-1})}$ to $\frac{1}{x_1(j_1)}, \dots, \frac{1}{x_{t-1}(j_{t-1})}$. We compute i_t, j_{t+1} , and a $\mu/2$ -approximation $\frac{1}{\tilde{x}_t(j_t)}$ to $\frac{1}{x_t(j_t)}$.

We first determine i_t in one pass. This can be done since A is presented in row order, together with reservoir sampling. Namely, given row A_k , we compute for each $1 \leq t' \leq t-1$, a μ -approximation $\tilde{v}_{t'}(k) = A_k(j_{t'}) \cdot \frac{1}{\tilde{x}_{t'}(j_{t'})}$ to $v_{t'}(k) = A_k(j_{t'}) \cdot \frac{1}{x_{t'}(j_{t'})}$, and then

$$\tilde{p}_t(k) = \frac{1}{n} \cdot \prod_{t'=1}^{t-1} (1 + \eta \tilde{v}_{t'}(k) + \eta^2 \tilde{v}_{t'}^2(k)).$$

Thus, we can reconstruct $\tilde{p}_t(k)$ for use with reservoir sampling to obtain a sample i_t .

In the next $O(\log n)$ passes we obtain j_{t+1} as follows. To ℓ_2 -sample from x_t , we use Theorem 5.2 to sample a coordinate from the length- $(t-1)d$ stream consisting of the entries of the concatenated list: $L = A_{i_1}, A_{i_2}, \dots, A_{i_{t-1}}$. Notice that $y_t = \frac{1}{\sqrt{2T}} \cdot \sum_{j=1}^{t-1} A_{i_j}$, and so Theorem 5.2 applied to L implements ℓ_2 -sampling from x_t . However, the algorithm returns $y_t(j_t)$ rather than $x_t(j_t)$. To obtain an approximation to $x_t(j_t)$, we $(\epsilon/3)$ -approximate $\|y_t\|$ using Theorem 5.1, from which $x_t(j_t) = \frac{y_t(j_t)}{\max\{1, \|y_t\|\}}$. We thus obtain an $(\epsilon/2)$ -approximation $\frac{1}{\tilde{x}_t(j_t)}$ to $\frac{1}{x_t(j_t)}$.

Using Lemma A.2, letting $y_{T+1} = \frac{1}{\sqrt{2T}} \sum_{j=1}^T A_{i_j}$, then $x_{T+1} = \frac{y_{T+1}}{\max\{1, \|y_{T+1}\|\}}$ results in an additive ϵ approximation. To compute this, we must $(1 \pm \epsilon)$ -approximate $\|y_{T+1}\|$, which we do in an additional pass using Theorem 5.1. Note that we cannot afford d space, which would be required to compute the norm exactly.

Theorem 5.3. *There is an $\tilde{O}(\epsilon^{-2})$ -pass, $\tilde{O}(\epsilon^{-2})$ -space algorithm running in total time $\tilde{O}(\epsilon^{-4}(n+d))$ which returns a list of $T = \tilde{O}(\epsilon^{-2})$ row indices i_1, \dots, i_T which implicitly represent the normal vector to a hyperplane for ϵ -approximate classification, together with an additive- ϵ approximation to the margin.*

For the MEB problem with high probability there are only $\tilde{O}(\epsilon^{-1})$ different values of i_t (i.e., updates to the primal vector). An important point is that we can get all $\tilde{O}(\epsilon^{-1})$ ℓ_2 -samples independently from the same primal vector between changes to it by running the algorithm of [MW10] independently $\tilde{O}(\epsilon^{-1})$ times in parallel.

We spend $\tilde{O}((n+d)\epsilon^{-2})$ time per iteration, to reconstruct the dual vector and run the algorithm of [MW10] independently $\tilde{O}(\epsilon^{-1})$ times on a stream of length $\tilde{O}(d\epsilon^{-1})$ to do ℓ_2 -sampling).

Minimum Enclosing Ball For the MEB problem we need the following standard tool.

Fact 5.4. *(see, e.g., [?]) Let $\sigma \in \{-1, 1\}^d$ be uniform from a 4-wise independent family of sign vectors. For any n -dimensional vector v , $\mathbf{E}_\sigma[\langle \sigma, v \rangle^2] = \|v\|_2^2$ and $\mathbf{Var}_\sigma[\langle \sigma, v \rangle^2] \leq 2\|v\|_2^4$.*

Define an epoch to be a contiguous block of iterations for which x_t does not change. Notice that x_t does not change with probability $1 - \alpha$.

We describe the necessary modifications to Algorithm 2. Throughout we round matrix entries to the nearest integer multiple of $\text{poly}(1/d)$ for a sufficiently large polynomial. We use the fact argued in §C that the algorithm remains correct given the per-iteration rounding of the updates $v_t(i)$ to relative error μ , where μ is on the order of $\eta\epsilon/T$.

We will not compute $\|x_t\|$ in each epoch. This would require $\Omega(d)$ space. However, unlike in the case of classification, for the MEB problem we cannot even afford to use Theorem 5.1 to approximate $\|x_t\|^2$, as that would cost $\Omega(\epsilon^{-2})$ space. Instead, we will use Fact 5.4 to obtain an

unbiased estimator of $\|x_t\|^2$, which suffices for our analysis to go through. Namely, by the triangle inequality, $\|x_t\| \leq 1$ (since we divide by t), and so the estimator of Fact 5.4 has variance $O(1)$.

Again, we implicitly represent the primal and dual vectors. We only store one index i_s and j_s per epoch s . In epoch s , we have indices i_1, \dots, i_s , which correspond to indices of the row A_i chosen for use to update the primal vector in the current and previous epochs. As in the non-streaming version of this algorithm, we use the same coordinate j_s for ℓ_2 -sampling in all iterations in an epoch and for all rows. Hence, throughout the course of the algorithm, the expected number of indices i_s and j_s that the algorithm stores is the number $\tilde{O}(\alpha T) = \tilde{O}(\varepsilon^{-1})$ of epochs. The algorithm also stores the number m_s of iterations in each epoch in the same amount of space.

At the beginning of the s -th epoch, we have maintained $\mu/2$ -approximations $\frac{1}{\bar{x}_1(j_1)}, \dots, \frac{1}{\bar{x}_{s-1}(j_{s-1})}$ to $\frac{1}{x_1(j_1)}, \dots, \frac{1}{x_{s-1}(j_{s-1})}$. We compute i_s, j_s , and a $\mu/2$ -approximation $\frac{1}{\bar{x}_s(j_s)}$ to $\frac{1}{x_s(j_s)}$.

We first determine i_s in one pass. This can be done as in classification since A is presented in row order, together with reservoir sampling. Namely, given row A_k , we compute for each $1 \leq s' \leq s-1$, a μ -approximation $\tilde{v}_{s'}(k) = A_k(j_{s'}) \cdot \frac{1}{\bar{x}_{s'}(j_{s'})}$ to $v_{s'}(k) = A_k(j_{s'}) \cdot \frac{1}{x_{s'}(j_{s'})}$, and then

$$\tilde{p}_s(k) = \frac{1}{n} \cdot \prod_{s'=1}^{s-1} (1 + \eta \tilde{v}_{s'}(k) + \eta^2 \tilde{v}_{s'}^2(k))^{m_{s'}}.$$

Thus, we can reconstruct $\tilde{p}_s(k)$ for use with reservoir sampling to obtain a sample i_s .

In the next $O(\log n)$ passes we obtain j_s as follows. To ℓ_2 -sample from x_s , we use Theorem 5.2 to sample a coordinate from the length- $(s-1)d$ stream consisting of the entries of the concatenated list: $L = A_{i_1}, A_{i_2}, \dots, A_{i_{s-1}}$. Notice that $y_s = \sum_{j=1}^{s-1} A_{i_j}$, and so Theorem 5.2 applied to L implements ℓ_2 -sampling from y_s , and hence x_s as well. We thus obtain an $(\varepsilon/2)$ -approximation $\frac{1}{\bar{x}_s(j_s)}$ to $\frac{1}{x_s(j_s)}$.

Applying Fact 5.4, we obtain

Theorem 5.5. *Given the norms of each row A_i , there is an $\tilde{O}(\varepsilon^{-1})$ -pass, $\tilde{O}(\varepsilon^{-2})$ -space algorithm running in total time $\tilde{O}(\varepsilon^{-3}(n+d))$ which returns a list of $T = \tilde{O}(\varepsilon^{-1})$ row indices i_1, \dots, i_T which implicitly represent the MEB center, together with an additive ε -approximation to the MEB radius.*

6 Kernelizing the Sublinear algorithms

An important generalization of linear classifiers is that of kernel-based linear predictors (see e.g. [SS03]). Let $\Psi : \mathbb{R}^d \mapsto \mathcal{H}$ be a mapping of feature vectors into a reproducing kernel Hilbert space. In this setting, we seek a non-linear classifier given by $h \in \mathcal{H}$ so as to maximize the margin:

$$\sigma \equiv \max_{h \in \mathcal{H}} \min_{i \in [n]} \langle h, \Psi(A_i) \rangle.$$

The kernels of interest are those for which we can compute inner products of the form $k(x, y) = \langle \Psi(x), \Psi(y) \rangle$ efficiently.

One popular kernel is the polynomial kernel, for which the corresponding Hilbert space is the set of polynomials over \mathbb{R}^d of degree q . The mapping Ψ for this kernel is given by

$$\forall S \subseteq [d], |S| \leq q. \Psi(x)_S = \prod_{i \in S} x_i.$$

That is, all monomials of degree at most q . The kernel function in this case is given by $k(x, y) = (x^\top y)^q$. Another useful kernel is the Gaussian kernel $k(x, y) = \exp(-\frac{\|x-y\|^2}{2s^2})$, where s is a parameter. The mapping here is defined by the kernel function (see [SS03] for more details).

Algorithm 4 Sublinear Kernel Perceptron

1: Input: $\varepsilon > 0$, $A \in \mathbb{R}^{n \times d}$ with $A_i \in \mathbb{B}$ for $i \in [n]$.
2: Let $T \leftarrow 200^2 \varepsilon^{-2} \log n$, $y_1 \leftarrow 0$, $w_1 \leftarrow \vec{1}_n$, $\eta \leftarrow \frac{1}{100} \sqrt{\frac{\log n}{T}}$.
3: **for** $t = 1$ to T **do**
4: $p_t \leftarrow \frac{w_t}{\|w_t\|_1}$, $x_t \leftarrow \frac{y_t}{\max\{1, \|y_t\|\}}$.
5: Choose $i_t \in [n]$ by $i_t \leftarrow i$ with probability $p_t(i)$.
6: $y_{t+1} \leftarrow \sum_{\tau \in [t]} \Psi(A_{i_\tau}) / \sqrt{2T}$.
7: **for** $i \in [n]$ **do**
8: $\tilde{v}_t(i) \leftarrow \mathbf{Kernel-L2-Sampling}(x_t, \Psi(A_i))$. (estimating $\langle x_t, \Psi(A_i) \rangle$)
9: $v_t(i) \leftarrow \text{clip}(\tilde{v}_t(i), 1/\eta)$.
10: $w_{t+1}(i) \leftarrow w_t(i)(1 - \eta v_t(i) + \eta^2 v_t(i)^2)$.
11: **end for**
12: **end for**
13: **return** $\bar{x} = \frac{1}{T} \sum_t x_t$

The kernel version of Algorithm 1 is shown in Figure 4. Note that x_t and y_t are members of \mathcal{H} , and not maintained explicitly, but rather are implicitly represented by the values i_t . (And thus $\|y_t\|$ is the norm of \mathcal{H} , not \mathbb{R}^d .) Also, $\Psi(A_i)$ is not computed. The needed kernel product $\langle x_t, \Psi(A_i) \rangle$ is estimated by the procedure **Kernel-L2-Sampling**, using the implicit representations and specific properties of the kernel being used. In the regular sublinear algorithm, this inner product could be sufficiently well approximated in $O(1)$ time via ℓ_2 -sampling. As we show below, for many interesting kernels the time for **Kernel-L2-Sampling** is not much longer.

For the analog of Theorem 2.7 to apply, we need the expectation of the estimates $v_t(i)$ to be correct, with variance $O(1)$. By Lemma C.1, it is enough if the estimates $v_t(i)$ have an additive bias of $O(\varepsilon)$. Hence, we define the procedure **Kernel-L2-Sampling** to obtain such an not-too-biased estimator with variance at most one; first we show how to implement **Kernel-L2-Sampling**, assuming that there is an estimator $\tilde{k}(\cdot)$ of the kernel $k(\cdot)$ such that $\mathbf{E}[\tilde{k}(x, y)] = k(x, y)$ and $\mathbf{Var}(\tilde{k}(x, y)) \leq 1$, and then we show how to implement such kernel estimators.

6.1 Implementing Kernel-L2-Sampling

Estimating $\|y\|_t$ A key step in **Kernel-L2-Sampling** is the estimation of $\|y_t\|$, which readily reduces to estimating

$$Y_t \equiv 2T \|y_t\|^2 / t^2 = \frac{1}{t^2} \sum_{\tau, \tau' \in [t]} k(A_{i_\tau}, A_{i_{\tau'}}),$$

that is, the mean of the summands. Since we use $\max\{1, \|y_t\|\}$, we need not be concerned with small $\|y_t\|$, and it is enough that the additive bias in our estimate of Y be at most $\varepsilon/T \leq \varepsilon(2T/t^2)$ for $t \in [T]$, implying a bias for $\|y_t\|$ no more than ε . Since we need $1/\|y_t\|$ in the algorithm, it is not enough for estimates of Y just to be good in mean and variance; we will find an estimator whose error bounds hold with high probability.

Our estimate \tilde{Y}_t of Y_t can first be considered assuming we only need to make an estimate for a single value of t .

Let $N_Y \leftarrow t^2 \lceil (8/3) \log(1/\delta) T^2 / \varepsilon^2 t^2 \rceil$. To estimate Y_t , we compute, for each $\tau, \tau' \in [t]$, $n_t \leftarrow N_Y / t^2$ independent estimates

$$X_{\tau, \tau', m} \leftarrow \text{clip}(\tilde{k}(A_{i_\tau}, A_{i_{\tau'}}), T/\varepsilon), \text{ for } m \in [n_t],$$

and our estimate is

$$\tilde{Y}_t \leftarrow \sum_{\substack{\tau, \tau' \in [t] \\ m \in [n_t]}} X_{\tau, \tau', m} / N_Y.$$

Lemma 6.1. *With probability at least $1 - \delta$, $|Y - \tilde{Y}_t| \leq \epsilon/T$.*

Proof. We apply Bernstein's inequality (as in 32) to the N_Y random variables $X_{\tau, \tau', m} - \mathbf{E}[X_{\tau, \tau', m}]$, which have mean zero, variance at most one, and are at most T/ϵ in magnitude. Bernstein's inequality implies, using $\text{Var}[X_{\tau, \tau', m}] \leq 1$,

$$\log \text{Prob}\left\{ \sum_{\substack{\tau, \tau' \in [t] \\ m \in [n_t]}} (X_{\tau, \tau', m} - \mathbf{E}[X_{\tau, \tau', m}]) > \alpha \right\} \leq -\alpha^2 / (N_Y + (T/\epsilon)\alpha/3),$$

and putting $\alpha \leftarrow N_Y \epsilon / T$ gives

$$\begin{aligned} \log \text{Prob}\{\tilde{Y} - \mathbf{E}[\tilde{Y}] > \epsilon/T\} &\leq -N_Y^2 (\epsilon/T)^2 / (N_Y + (T/\epsilon)N_Y (\epsilon/T)/3) \\ &\leq -(8/3) \log(1/\delta) (3/4) \leq -2 \log(1/\delta). \end{aligned}$$

Similar reasoning for $-X_{\tau, \tau', m}$, and the union bound, implies the lemma. \square

To compute Y for $t = 1 \dots T$, we can save some work by reusing estimates from one t to the next. Now let $N_Y \leftarrow \lceil (8/3) \log(1/\delta) T^2 / \epsilon^2 \rceil$. Compute \tilde{Y}_1 as above for $t = 1$, and let $\hat{Y}_1 \leftarrow \tilde{Y}_1$. For $t > 1$, let $n_t \leftarrow \lceil N_Y / t^2 \rceil$, and let

$$\hat{Y}_t \leftarrow \sum_{m \in [n_t]} X_{t, t, m} / n_t + \sum_{\substack{\tau \in [t] \\ m \in [n_t]}} (X_{t, \tau, m} + X_{\tau, t, m}) / n_t,$$

and return $\tilde{Y}_t \leftarrow \sum_{\tau \in [t]} \hat{Y}_\tau / t^2$.

Since for each τ and τ' , the expected total contribution of all $X_{\tau, \tau', m}$ terms to \tilde{Y}_t is $k(A_{i_\tau}, A_{i_{\tau'}})$, we have $\mathbf{E}[\tilde{Y}_t] = Y_t$. Moreover, the number of instances of $X_{\tau, \tau', m}$ averaged to compute \tilde{Y}_t is always at least as large as the number used for the above ‘‘batch’’ version; it follows that the total variance of \tilde{Y}_t is non-increasing in t , and therefore Lemma 6.1 holds also for the \tilde{Y}_t computed stepwise.

Since the number of calls to $\tilde{k}(\cdot)$ is $\sum_{t \in [T]} (1 + 2n_t) = O(N_Y)$, we have the following lemma.

Lemma 6.2. *The values $\tilde{Y}_t(t^2/2T) \approx \|y_t\|$, $t \in [T]$, can be estimated with $O((\log(1/\epsilon\delta)T^2/\epsilon^2))$ calls to $\tilde{k}(\cdot)$, so that with probability at least $1 - \delta$, $|\tilde{Y}_t(t^2/2T) - \|y_t\|| \leq \epsilon$. The values $\|y_t\|$, $t \in [T]$, can be computed exactly with T^2 calls to the exact kernel $k(\cdot, \cdot)$.*

Proof. This follows from the discussion above, applying the union bound over $t \in [T]$, and adjusting constants. The claim for exact computation is straightforward. \square

Given this procedure for estimating $\|y_t\|$, we can describe **Kernel-L2-Sampling**. Since $x_{t+1} = y_{t+1} / \max\{1, \|y_{t+1}\|\}$, we have

$$\begin{aligned} \langle x_{t+1}, A_i \rangle &= \frac{1}{\max\{1, \|y_{t+1}\|\} \sqrt{2T}} \sum_{\tau \in [t]} \langle \Psi(A_{i_\tau}), \Psi(A_i) \rangle \\ &= \frac{1}{\max\{1, \|y_{t+1}\|\} \sqrt{2T}} \sum_{\tau \in [t]} k(A_{i_\tau}, A_i), \end{aligned} \tag{23}$$

so that the main remaining step is to estimate $\sum_{\tau \in [t]} k(A_{i_\tau}, A_i)$, for $i \in [n]$. Here we simply call $\tilde{k}(A_{i_\tau}, A_i)$ for each τ . We save time, at the cost of $O(n)$ space, by saving the value of the sum for each $i \in [n]$, and updating it for the next t with n calls $\tilde{k}(A_{i_t}, A_i)$.

Lemma 6.3. *Let L_k denote the expected time needed for one call to $\tilde{k}(\cdot, \cdot)$, and T_k denote the time needed for one call to $k(\cdot, \cdot)$. Except for estimating $\|y_t\|$, **Kernel-L2-Sampling** can be computed in nL_k expected time per iteration t . The resulting estimate has expectation within additive ϵ of $\langle x_t, A_i \rangle$, and variance at most one. Thus Algorithm 4 runs in time $\tilde{O}(\frac{(L_k n + d)}{\epsilon^2} + \min\{\frac{L_k}{\epsilon^6}, \frac{T_k}{\epsilon^4}\})$, and produces a solution with properties as in Algorithm 1.*

Proof. For **Kernel-L2-Sampling** it remains only to show that its variance is at most one, given that each $\tilde{k}(\cdot, \cdot)$ has variance at most one. We observe from (23) that t independent estimates $\tilde{k}(\cdot, \cdot)$ are added together, and scaled by a value that is at most $1/\sqrt{2T}$. Since the variance of the sum is at most t , and the variance is scaled by a value no more than $1/2T$, the variance of **Kernel-L2-Sampling** is at most one. The only bias in the estimate is due to estimation of $\|y_t\|$, which gives relative error of ϵ . For our kernels, $\|\Psi(v)\| \leq 1$ if $v \in \mathbb{B}$, so the additive error of **Kernel-L2-Sampling** is $O(\epsilon)$.

The analysis of Algorithm 4 then follows as for the un-kernelized perceptron; we neglect the time needed for preprocessing for the calls to $\tilde{k}(\cdot, \cdot)$, as it is dominated by other terms for the kernels we consider, and this is likely in general. \square

6.2 Implementing the Kernel Estimators

Using the lemma above we can derive corollaries for the Gaussian and polynomial kernels.

Polynomial kernels For the polynomial kernel of degree q , estimating a single kernel product, i.e. $k(x, y) = k(A_i, A_j)$, where the norm of x, y is at most one, takes $O(q)$ as follows: Recall that for the polynomial kernel, $k(x, y) = (x^\top y)^q$. To estimate this kernel we take the product of q independent ℓ_2 -samples, yielding $\tilde{k}(x, y)$. Notice that the expectation of this estimator is exactly equal to the product of expectations, $\mathbf{E}[\tilde{k}(x, y)] = (x^\top y)^q$. The variance of this estimator is equal to the product of variances, which is $\mathbf{Var}(\tilde{k}(x, y)) \leq (\|x\| \|y\|)^q \leq 1$. Of course, calculating the inner product exactly takes $O(d \log q)$ time. We obtain:

Corollary 6.4. *For the polynomial degree- q kernel, Algorithm 4 runs in time*

$$\tilde{O}\left(\frac{q(n+d)}{\epsilon^2} + \min\left\{\frac{d \log q}{\epsilon^4}, \frac{q}{\epsilon^6}\right\}\right).$$

Gaussian kernels To estimate the Gaussian kernel function, we assume that $\|x\|$ and $\|y\|$ are known and no more than $s/2$; thus to estimate

$$k(x, y) = \exp(-\|x - y\|^2) = \exp(-(\|x\|^2 + \|y\|^2)/2s^2) \exp(x^\top y/s^2),$$

we need to estimate $\exp(x^\top y/s^2)$. For $\exp(\gamma X) = \sum_{i \geq 0} \gamma^i X^i / i!$ with random X and parameter $\gamma > 0$, we pick index i with probability $\exp(-\gamma) \gamma^i / i!$ (that is, i has a Poisson distribution) and return $\exp(\gamma)$ times the product of i independent estimates of X .

In our case we take X to be the average of c ℓ_2 -samples of $x^\top y$, and hence $\mathbf{E}[X] = x^\top y$, $\mathbf{E}[X^2] \leq \frac{1}{c} \mathbf{E}[(x^\top y)^2] \leq \frac{1}{c}$. The expectation of our kernel estimator is thus:

$$\mathbf{E}[\tilde{k}(x, y)] = \mathbf{E}\left[\sum_{i \geq 0} e^{-\gamma} \gamma^i i! \cdot e^\gamma \cdot X^i\right] = \sum_{i \geq 0} \gamma^i i! \prod_{j=1}^i \mathbf{E}[X] = \exp(\gamma x^\top y).$$

The second moment of this estimator is bounded by:

$$\mathbf{E}[\tilde{k}(x, y)^2] = \mathbf{E}\left[\sum_{i \geq 0} e^{-\gamma} \gamma^i i! \cdot e^{2\gamma} \cdot (X^i)^2\right] = e^\gamma \sum_{i \geq 0} \gamma^i i! \prod_{j=1}^i \mathbf{E}[X^2] \leq \exp\left(\frac{2\gamma}{c}\right).$$

Hence, we take $\gamma = c = \frac{1}{s^2}$. This gives a correct estimator in terms of expectation and constant variance. The variance can be further made smaller than one by taking the average of a constant estimators of the above type.

As for evaluation time, the expected size of the index i is $\gamma = \frac{1}{s^2}$. Thus, we require on the expectation $\gamma \times c = \frac{1}{s^4}$ of ℓ_2 -samples.

We obtain:

Corollary 6.5. *For the Gaussian kernel with parameter s , Algorithm 4 runs in time*

$$\tilde{O}\left(\frac{(n+d)}{s^4 \varepsilon^2} + \min\left\{\frac{d}{\varepsilon^4}, \frac{1}{s^4 \varepsilon^6}\right\}\right).$$

6.3 Kernelizing the MEB and strictly convex problems

Analogously to Algorithm 4, we can define the kernel version of strongly convex problems, including MEB. The kernelized version of MEB is particularly efficient, since as in Algorithm 2, the norm $\|y_t\|$ is never required. This means that the procedure **Kernel-L2-Sampling** can be computed in time $O(nL_k)$ per iteration, for a total running time of $O(L_k(\varepsilon^{-2}n + \varepsilon^{-1}d))$.

7 Lower bounds

All of our lower bounds are information-theoretic, meaning that any successful algorithm must read at least some number of entries of the input matrix A . Clearly this also lower bounds the time complexity of the algorithm in the unit-cost RAM model.

Some of our arguments use the following meta-theorem. Consider a $p \times q$ matrix A , where p is an even integer. Consider the following random process. Let $W \geq q$. Let $a = 1 - 1/W$, and let e_j denote the j -th standard q -dimensional unit vector. For each $i \in [p/2]$, choose a random $j \in [q]$ uniformly, and set $A_{i+p/2} \leftarrow A_i \leftarrow ae_j + b(\mathbf{1}_q - e_j)$, where b is chosen so that $\|A_i\|_2 = 1$. We say that such an A is a YES instance. With probability $1/2$, transform A into a NO instance as follows: choose a random $i^* \in [p/2]$ uniformly, and if $A_{i^*} = ae_{j^*} + b(\mathbf{1}_q - e_{j^*})$ for a particular $j^* \in [q]$, set $A_{i^*+p/2} \leftarrow -ae_{j^*} + b(\mathbf{1}_q - e_{j^*})$.

Suppose there is a randomized algorithm reading at most s positions of A which distinguishes YES and NO instances with probability $\geq 2/3$, where the probability is over the algorithm's coin tosses and this distribution μ on YES and NO instances. By averaging this implies a deterministic algorithm Alg reading at most s positions of A and distinguishing YES and NO instances with probability $\geq 2/3$, where the probability is taken only over μ . We show the following meta-theorem with a standard argument.

Theorem 7.1. (Meta-theorem) *For any such algorithm Alg , $s = \Omega(pq)$.*

This Meta-Theorem follows from the following folklore fact:

Fact 7.2. *Consider the following random process. Initialize a length- r array A to an array of r zeros. With probability $1/2$, choose a random position $i \in [r]$ and set $A[i] = 1$. With the remaining probability $1/2$, leave A as the all zero array. Then any algorithm which determines if A is the all zero array with probability $\geq 2/3$ must read $\Omega(r)$ entries of A .*

Let us prove Theorem 7.1 using this fact:

Proof. Consider the matrix $B \in \mathbb{R}^{(p/2) \times q}$ which is defined by subtracting the “bottom” half of the matrix from the top half, that is, $B_{i,j} = A_{i,j} - A_{i+p/2,j}$. Then B is the all zeros matrix, except that with probability 1/2, there is one entry whose value is roughly two, and whose location is random and distributed uniformly. An algorithm distinguishing between YES and NO instances of A in particular distinguishes between the two cases for B , which cannot be done without reading a linear number of entries. \square

In the proofs of Theorem 7.3, Corollary 7.4, and Theorem 7.6, it will be more convenient to use M as an upper bound on the number of non-zero entries of A rather than the exact number of non-zero entries. However, it should be understood that these theorems (and corollary) hold even when M is exactly the number of non-zero entries of A .

To see this, our random matrices A constructed in the proofs have at most M non-zero entries. If this number M' is strictly less than M , we arbitrarily replace $M - M'$ zero entries with the value $(nd)^{-C}$ for a large enough constant $C > 0$. Under our assumptions on the margin or the minimum enclosing ball radius of the points, the solution value changes by at most a factor of $(1 \pm (nd)^{1-C})$, which does not affect the proofs.

7.1 Classification

Recall that the margin $\sigma(A)$ of an $n \times d$ matrix A is given by $\max_{x \in \mathbb{B}} \min_i A_i x$. Since we assume that $\|A_i\|_2 \leq 1$ for all i , we have that $\sigma(A) \leq 1$.

7.1.1 Relative Error

We start with a theorem for relative error algorithms.

Theorem 7.3. *Let $\kappa > 0$ be a sufficiently small constant. Let ε and $\sigma(A)$ have $\sigma(A)^{-2}\varepsilon^{-1} \leq \kappa \min(n, d)$, $\sigma(A) \leq 1 - \varepsilon$, with ε also bounded above by a sufficiently small constant. Also assume that $M \geq 2(n + d)$, that $n \geq 2$, and that $d \geq 3$. Then any randomized algorithm which, with probability at least 2/3, outputs a number in the interval $[\sigma(A) - \varepsilon\sigma(A), \sigma(A)]$ must read*

$$\Omega(\min(M, \sigma(A)^{-2}\varepsilon^{-1}(n + d)))$$

entries of A . This holds even if $\|A_i\|_2 = 1$ for all rows A_i .

Notice that this yields a stronger theorem than assuming that both n and d are sufficiently large, since one of these values may be constant.

Proof. We divide the analysis into cases: the case in which d or n is constant, and the case in which each is sufficiently large. Let $\tau \in [0, 1 - \varepsilon]$ be a real number to be determined.

Case: d or n is a constant By our assumption that $\sigma(A)^{-2}\varepsilon^{-1} \leq \kappa \min(n, d)$, the values $\sigma(A)$ and ε are constant, and sufficiently large. Therefore we just need to show an $\Omega(\min(M, n + d))$ bound on the number of entries read. By the premise of the theorem, $M = \Omega(n + d)$, so we can just show an $\Omega(n + d)$ bound.

An $\Omega(d)$ bound. We give a randomized construction of an $n \times d$ matrix A .

The first row of A is built as follows. Let $A_{1,1} \leftarrow \tau$ and $A_{1,2} \leftarrow 0$. Pick $j^* \in \{3, 4, \dots, d\}$ uniformly at random, and let $A_{1,j^*} \leftarrow \varepsilon^{1/2}\tau$. For all remaining $j \in \{3, 4, \dots, d\}$, assign $A_{1,j} \leftarrow \zeta$,

where $\zeta \leftarrow 1/d^3$. (The role of ζ is to make an entry slightly non-zero to prevent an algorithm which has access to exactly the non-zero entries from skipping over it.) Now using the conditions on τ , we have

$$X \leftarrow \|A_1\|^2 = \tau^2 + (d-3)\zeta^2 + \varepsilon\tau^2 \leq (1-\varepsilon)^2 + d^{-2} + \varepsilon \leq 1 - \varepsilon + \varepsilon^2 + \kappa^2\varepsilon^2 \leq 1,$$

and so by letting $A_{1,2} \leftarrow \sqrt{1-X}$, we have $\|A_1\| = 1$.

Now we let $A_2 \leftarrow -A_1$, with two exceptions: we let $A_{2,1} \leftarrow A_{1,1} = \tau$, and with probability $1/2$, we negate A_{2,j^*} . Thus $\|A_2\| = 1$ also.

For row i with $i > 2$, put $A_{i,1} \leftarrow (1+\varepsilon)\tau$, $A_{i,2} \leftarrow \sqrt{1-A_{i,1}^2}$, and all remaining entries zero.

We have the following picture.

$$\begin{pmatrix} \tau & (1-\tau^2-(d-3)\zeta^2-\varepsilon\tau^2)^{1/2} & \zeta & \cdots & \zeta & \varepsilon^{1/2}\tau & \zeta & \cdots & \zeta \\ \tau & -(1-\tau^2-(d-3)\zeta^2-\varepsilon\tau^2)^{1/2} & -\zeta & \cdots & -\zeta & \pm\varepsilon^{1/2}\tau & -\zeta & \cdots & -\zeta \\ (1+\varepsilon)\tau & (1-(1+\varepsilon)^2\tau^2)^{1/2} & 0 & \cdots & & & & & 0 \\ (1+\varepsilon)\tau & (1-(1+\varepsilon)^2\tau^2)^{1/2} & 0 & \cdots & & & & & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & & & & \cdots \\ (1+\varepsilon)\tau & (1-(1+\varepsilon)^2\tau^2)^{1/2} & 0 & \cdots & & & & & 0 \end{pmatrix}$$

Observe that the the number of non-zero entries of the resulting matrix is $2n + 2d - 4$, which satisfies the premise of the theorem. Moreover, all rows A_i satisfy $\|A_i\| = 1$.

Notice that if $A_{1,j^*} = -A_{2,j^*}$, then the margin of A is at most τ , which follows by observing that all but the first coordinate of A_1 and A_2 have opposite signs.

On the other hand, if $A_{1,j^*} = A_{2,j^*}$, consider the vector y with $y_1 \leftarrow 1$, $y_{j^*} \leftarrow \sqrt{\varepsilon}$, and all other entries zero. Then for all i , $A_i y = \tau(1+\varepsilon)$, and so the unit vector $x \leftarrow y/\|y\|$ has

$$A_i x = \frac{\tau(1+\varepsilon)}{\sqrt{1+\varepsilon}} = \tau(1+\varepsilon)^{1/2} = \tau(1+\Omega(\varepsilon)).$$

It follows that in this case the margin of A is at least $\tau(1+\Omega(\varepsilon))$. Setting $\tau = \Theta(\sigma)$ and rescaling ε by a constant factor, it follows that these two cases can be distinguished by an algorithm satisfying the premise of the theorem. By Fact 7.2, any algorithm distinguishing these two cases with probability $\geq 2/3$ must read $\Omega(d)$ entries of A .

An $\Omega(n)$ bound. We construct the $n \times d$ matrix A as follows. All but the first two columns are 0. We set $A_{i,1} \leftarrow \tau$ and $A_{i,2} \leftarrow \sqrt{1-\tau^2}$ for all $i \in [n]$. Next, with probability $1/2$, we pick a random row i^* , and negate $A_{i^*,2}$. We have the following picture.

$$\begin{pmatrix} \tau & \sqrt{1-\tau^2} & 0 & \cdots & 0 \\ \cdots & \cdots & 0 & \cdots & 0 \\ \tau & \sqrt{1-\tau^2} & 0 & \cdots & 0 \\ \tau & \pm\sqrt{1-\tau^2} & 0 & \cdots & 0 \\ \tau & \sqrt{1-\tau^2} & 0 & \cdots & 0 \\ \cdots & \cdots & 0 & \cdots & 0 \\ \tau & \sqrt{1-\tau^2} & 0 & \cdots & 0 \end{pmatrix}$$

The number of non-zeros of the resulting matrix is $2n < M$. Depending on the sign of $A_{i^*,2}$, the margin of A is either 1 or τ . Setting $\tau = \Theta(\sigma)$, an algorithm satisfying the premise of the theorem can distinguish the two cases. By Fact 7.2, any algorithm distinguishing these two cases with probability $\geq 2/3$ must read $\Omega(n)$ entries of A .

Case: d and n are both sufficiently large Suppose first that $M = \Omega(\sigma(A)^{-2}\varepsilon^{-1}(n+d))$ for a sufficiently large constant in the $\Omega(\cdot)$. Let s be an even integer in $\Theta(\tau^{-2}\varepsilon^{-1})$ and with $s < \min(n, d) - 1$. We will also choose a value τ in $\Theta(\sigma(A))$. We can assume without loss of generality that n and d are sufficiently large, and even.

An $\Omega(ns)$ bound. We set the d -th entry of each row of A to the value τ . We set all entries in columns $s+1$ through $d-1$ to 0. We then choose the remaining entries of A as follows. We apply Theorem 7.1 with parameters $p = n, q = s$, and $W = d^2$, obtaining an $n \times s$ matrix B , where $\|B_i\| = 1$ for all rows B_i . Put $B' \leftarrow B\sqrt{1-\tau^2}$. We then set $A_{i,j} \leftarrow B'_{i,j}$ for all $i \in [n]$ and $j \in [s]$. We have the following block structure for A .

$$\begin{bmatrix} B\sqrt{1-\tau^2} & \mathbf{0}_{n \times (d-s-1)} & \mathbf{1}_n\tau \end{bmatrix}$$

Here $\mathbf{0}_{n \times (d-s-1)}$ is a matrix of all 0's, of the given dimensions. Notice that $\|A_i\| = 1$ for all rows A_i , and the number of non-zero entries is at most $n(s+1)$, which is less than the value M .

We claim that if B is a YES instance, then the margin of A is $\tau(1 + \Omega(\varepsilon))$. Indeed, consider the unit vector x for which

$$x_j \leftarrow \begin{cases} \left(\frac{\varepsilon}{s} - \frac{\varepsilon^2}{4s}\right)^{1/2} & j \in [s] \\ 0 & j \in [s+1, d-1] \\ 1 - \varepsilon/2 & j = d \end{cases} \quad (24)$$

For any row A_i ,

$$\begin{aligned} A_i x &\geq \left(\frac{\varepsilon}{s} - \frac{\varepsilon^2}{4s}\right)^{1/2} \left(\sqrt{1-\tau^2} - O\left(\frac{\sqrt{1-\tau^2}}{d^2}\right)\right) + \left(1 - \frac{\varepsilon}{2}\right)\tau \\ &\geq \left(\frac{\varepsilon}{s} - \frac{\varepsilon^2}{4s}\right)^{1/2} \left(1 - \tau - O\left(\frac{\sqrt{1-\tau^2}}{d^2}\right)\right) + \tau - \frac{\varepsilon\tau}{2} && \text{since } \sqrt{1-\tau^2} \geq 1 - \tau \\ &\geq \left(\frac{\varepsilon}{s}\right)^{1/2} (1 - \tau) + \tau - \frac{\varepsilon\tau}{2} - O(\varepsilon^2\tau^2) && \text{since } \sqrt{\frac{\varepsilon}{s}} \cdot \frac{1}{d^2} = O(\varepsilon^2\tau^2) \end{aligned}$$

If we set $s = c\tau^{-2}\varepsilon^{-1}$ for $c \in (0, 4)$, then

$$A_i x \geq \tau + \frac{\tau\varepsilon}{c^{1/2}} - \tau \left(\frac{\varepsilon}{2} + \frac{\tau\varepsilon}{c^{1/2}}\right) - O(\varepsilon^2\tau^2) = \tau(1 + \Omega(\varepsilon)). \quad (25)$$

On the other hand, if B is a NO instance, we claim that the margin of A is at most $\tau(1 + O(\varepsilon^2))$. By definition of a NO instance, there are rows A_i and A_j of A which agree except on a single column k , for which $A_{i,k} = \sqrt{1-\tau^2} - O\left(\frac{1-\tau^2}{d^2}\right)$ while $A_{j,k} = -A_{i,k}$. It follows that the x which maximizes $\min\{A_i x, A_j x\}$ has $x_k = 0$. But $\sum_{k' \neq k} A_{i,k'}^2 = 1 - (1 - \tau^2) + O\left(\frac{1}{d^2}\right) = \tau^2 + O\left(\frac{1}{d^2}\right)$. Since $\|x\| \leq 1$, by the Cauchy-Schwarz inequality

$$A_i x = A_j x \leq \left(\tau^2 + O\left(\frac{1}{d^2}\right)\right)^{1/2} \leq \tau + O(\varepsilon^2) = \tau(1 + O(\varepsilon^2)), \quad (26)$$

where the first inequality follows from our bound $\tau^{-2}\varepsilon^{-1} = O(d)$.

Setting $\tau = \Theta(\sigma(A))$ and rescaling ε by a constant factor, an algorithm satisfying the premise of the theorem can distinguish the two cases, and so by Theorem 7.1, it must read $\Omega(ns) =$

$\Omega(\sigma(A)^{-2}\varepsilon^{-1}n)$ entries of A .

An $\Omega(ds)$ bound. We first define rows $s + 1$ through n of our $n \times d$ input matrix A . For $i > s$, put $A_{i,d} \leftarrow \tau(1 + \varepsilon)$, $A_{i,d-1} \leftarrow (1 - \tau^2(1 + \varepsilon)^2)^{1/2}$, and all remaining entries zero.

We now define rows 1 through s . Put $A_{i,d} \leftarrow \tau$ for all $i \in [s]$. Now we apply Theorem 7.1 with $p = s$, $q = d - 2$, and $W = d^2$, obtaining an $s \times (d - 2)$ matrix B , where $\|B_i\| = 1$ for all rows B_i . Put $B' \leftarrow B\sqrt{1 - \tau^2}$, and set $A_{i,j} \leftarrow B'_{i,j}$ for all $i \in [s]$ and $j \in [d - 2]$. We have the following block structure for A .

$$\begin{bmatrix} B\sqrt{1 - \tau^2} & \mathbf{0}_s & \mathbf{1}_s\tau \\ \mathbf{0}_{(n-s) \times (d-2)} & \mathbf{1}_{n-s}(1 - \tau^2(1 + \varepsilon)^2)^{1/2} & \mathbf{1}_{n-s}\tau(1 + \varepsilon) \end{bmatrix}$$

Notice that $\|A_i\| = 1$ for all rows A_i , and the number of non-zero entries is at most $2n + sd < M$.

If B is a YES instance, let x be as in Equation (24). Since the first s rows of A agree with those in our proof of the $\Omega(ns)$ bound, then as shown in Equation (25), $A_i x = \tau(1 + \Omega(\varepsilon))$ for $i \in [s]$. Moreover, for $i > s$, since YES instances B are entry-wise positive, we have

$$A_i x > \left(1 - \frac{\varepsilon}{2}\right) \cdot \tau(1 + \varepsilon) = \tau(1 + \Omega(\varepsilon)).$$

Hence, if B is a YES instance the margin is $\tau(1 + \Omega(\varepsilon))$.

Now suppose B is a NO instance. Then, as shown in Equation (26), for any x for which $\|x\| \leq 1$, we have $A_i x \leq \tau(1 + O(\varepsilon^2))$ for $i \in [s]$. Hence, if B is a NO instance, the margin is at most $\tau(1 + O(\varepsilon^2))$.

Setting $\tau = \Theta(\sigma(A))$ and rescaling ε by a constant factor, an algorithm satisfying the premise of the theorem can distinguish the two cases, and so by Theorem 7.1, it must read $\Omega(ds) = \Omega(\sigma(A)^{-2}\varepsilon^{-1}d)$ entries of A .

Finally, if $M = O((n + d)\sigma(A)^{-2}\varepsilon^{-1})$, then we must show an $\Omega(M)$ bound. We will use our previous construction for showing an $\Omega(ns)$ bound, but replace the value of n there with n' , where n' is the largest integer for which $n's \leq M/2$. We claim that $n' \geq 1$. To see this, by the premise of the theorem $M \geq 2(n + d)$. Moreover, $s = \Theta(\varepsilon^{-1})$ and $\varepsilon^{-1} \leq \kappa(n + d)$. For a small enough constant $\kappa > 0$, $s \leq (n + d) \leq M/2$, as needed.

As the theorem statement concerns matrices with n rows, each of unit norm, we must have an input A with n rows. To achieve this, we put $A_{i,d} = \tau(1 + \varepsilon)$ and $A_{i,d-1} = (1 - \tau^2(1 + \varepsilon)^2)^{1/2}$ for all $i > n'$. In all remaining entries in rows A_i with $i > n'$, we put the value 0. This ensures that $\|A_i\| = 1$ for all $i > n'$, and it is easy to verify that this does not change the margin of A . Hence, the lower bound is $\Omega(n's) = \Omega(M)$. Notice that the number of non-zero entries is at most $2n + n's \leq 2M/3 + M/3 = M$, as needed.

This completes the proof. \square

7.1.2 Additive Error

Here we give a lower bound for the additive error case. We give two different bounds, one when $\varepsilon < \sigma$, and one when $\varepsilon \geq \sigma$. Notice that $\sigma \geq 0$ since we may take the solution $x = \mathbf{0}_d$. The following is a corollary of Theorem 7.3.

Corollary 7.4. *Let $\kappa > 0$ be a sufficiently small constant. Let $\varepsilon, \sigma(A)$ be such that $\sigma(A)^{-1}\varepsilon^{-1} \leq \kappa \min(n, d)$ and $\sigma(A) \leq 1 - \varepsilon/\sigma(A)$, where $0 < \varepsilon \leq \kappa'\sigma$ for a sufficiently small constant $\kappa' > 0$. Also assume that $M \geq 2(n + d)$, $n \geq 2$, and $d \geq 3$. Then any randomized algorithm which, with probability at least $2/3$, outputs a number in the interval $[\sigma - \varepsilon, \sigma]$ must read*

$$\Omega(\min(M, \sigma^{-1}\varepsilon^{-1}(n + d)))$$

entries of A . This holds even if $\|A_i\| = 1$ for all rows A_i .

Proof. We simply set the value of ε in Theorem 7.3 to ε/σ . Notice that ε is at most a sufficiently small constant and the value $\sigma^{-2}\varepsilon^{-1}$ in Theorem 7.3 equals $\sigma^{-1}\varepsilon^{-1}$, which is at most $\kappa \min(n, d)$ by the premise of the corollary, as needed to apply Theorem 7.3. \square

The following handles the case when $\varepsilon = \Omega(\sigma)$.

Corollary 7.5. *Let $\kappa > 0$ be a sufficiently small constant. Let $\varepsilon, \sigma(A)$ be such that $\varepsilon^{-2} \leq \kappa \min(n, d)$, $\sigma(A) + \varepsilon < \frac{1}{\sqrt{2}}$, and $\varepsilon = \Omega(\sigma)$. Also assume that $M \geq 2(n + d)$, $n \geq 2$, and $d \geq 3$. Then any randomized algorithm which, with probability at least $2/3$, outputs a number in the interval $[\sigma - \varepsilon, \sigma]$ must read*

$$\Omega(\min(M, \varepsilon^{-2}(n + d)))$$

entries of A . This holds even if $\|A_i\| = 1$ for all rows A_i .

Proof. The proof is very similar to that of Theorem 7.3, so we just outline the differences. In the case that d or n is constant, we have the following families of hard instances:

An $\Omega(n)$ bound for constant d :

$$\begin{pmatrix} \tau & (1 - \tau^2 - (d - 3)\zeta^2 - 2(\varepsilon + \tau)^2)^{1/2} & \zeta & \cdots & \zeta & \sqrt{2}(\varepsilon + \tau) & \zeta & \cdots & \zeta \\ \tau & -(1 - \tau^2 - (d - 3)\zeta^2 - 2(\varepsilon + \tau)^2)^{1/2} & -\zeta & \cdots & -\zeta & \pm\sqrt{2}(\varepsilon + \tau) & -\zeta & \cdots & -\zeta \\ \sqrt{2}(\varepsilon + \tau) & (1 - 2(\varepsilon + \tau)^2)^{1/2} & 0 & \cdots & & & & & 0 \\ \sqrt{2}(\varepsilon + \tau) & (1 - 2(\varepsilon + \tau)^2)^{1/2} & 0 & \cdots & & & & & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & & & & \cdots \\ \sqrt{2}(\varepsilon + \tau) & (1 - 2(\varepsilon + \tau)^2)^{1/2} & 0 & \cdots & & & & & 0 \end{pmatrix}$$

An $\Omega(d)$ bound for constant n :

$$\begin{pmatrix} \tau & \sqrt{1 - \tau^2} & 0 & \cdots & 0 \\ \cdots & \cdots & 0 & \cdots & 0 \\ \tau & \sqrt{1 - \tau^2} & 0 & \cdots & 0 \\ \tau & \pm\sqrt{1 - \tau^2} & 0 & \cdots & 0 \\ \tau & \sqrt{1 - \tau^2} & 0 & \cdots & 0 \\ \cdots & \cdots & 0 & \cdots & 0 \\ \tau & \sqrt{1 - \tau^2} & 0 & \cdots & 0 \end{pmatrix}$$

In these two cases, depending on the sign of the undetermined entry the margin is either τ or at least $\tau + \varepsilon$ (in the $\Omega(d)$ bound, it is τ or 1, but we assume $\tau + \varepsilon < \frac{1}{\sqrt{2}}$). It follows for $\tau = \sigma(A)$, the algorithm of the corollary can distinguish these two cases, for which the lower bounds follow from the proof of Theorem 7.3.

For the case of n and d sufficiently large, we have the following families of hard instances. In each case, the matrix B is obtained by invoking Theorem 7.1 with the value of $s = \Theta(\varepsilon^{-2})$.

An $\Omega(n\varepsilon^{-2})$ bound for n, d sufficiently large:

$$\begin{bmatrix} B\sqrt{1 - \tau^2} & \mathbf{0}_{n \times (d-s-1)} & \mathbf{1}_n \tau \end{bmatrix}$$

An $\Omega(d\varepsilon^{-2})$ bound for n, d sufficiently large:

$$\begin{bmatrix} B\sqrt{1-\tau^2} & \mathbf{0}_s & \mathbf{1}_s\tau \\ \mathbf{0}_{(n-s)\times(d-2)} & \mathbf{1}_{n-s}(1-(\tau+\varepsilon)^2)^{1/2} & \mathbf{1}_{n-s}(\tau+\varepsilon) \end{bmatrix}$$

In these two cases, by setting $W = \text{poly}(nd)$ to be sufficiently large in Theorem 7.1, depending on whether B is YES or a NO instance the margin is either at most $\tau + \frac{1}{\text{poly}(nd)}$ or at least $\tau + \sqrt{1-\tau^2} \cdot 2\varepsilon$ (for an appropriate choice of s). For $\tau < 1/\sqrt{2}$, the algorithm of the corollary can distinguish these two cases, and therefore needs $\Omega(ns)$ time in the first case, and $\Omega(ds)$ time in the second.

The extension of the proofs to handle the case $M = o((n+d)\varepsilon^{-2})$ is identical to that given in the proof of Theorem 7.3. \square

7.2 Minimum Enclosing Ball

We start by proving the following lower bound for estimating the squared MEB radius to within an additive ε . In the next subsection we improve the $\Omega(\varepsilon^{-1}n)$ term in the lower bound to $\tilde{\Omega}(\varepsilon^{-2}n)$ for algorithms that either additionally output a coreset, or output a MEB center that is a convex combination of the input points. As our primal-dual algorithm actually outputs a coreset, as well as a MEB center that is a convex combination of the input points, those bounds apply to it. Our algorithm has both of these properties though satisfying one or the other would be enough to apply the lower bound. Together with the $\varepsilon^{-1}d$ bound given by the next theorem, these bounds establish its optimality.

Theorem 7.6. *Let $\kappa > 0$ be a sufficiently small constant. Assume $\varepsilon^{-1} \leq \kappa \min(n, d)$ and ε is less than a sufficiently small constant. Also assume that $M \geq 2(n+d)$ and that $n \geq 2$. Then any randomized algorithm which, with probability at least $2/3$, outputs a number in the interval*

$$\left[\min_x \max_i \|x - A_i\|^2 - \varepsilon, \min_x \max_i \|x - A_i\|^2 \right]$$

must read

$$\Omega(\min(M, \varepsilon^{-1}(n+d)))$$

entries of A . This holds even if $\|A_i\| = 1$ for all rows A_i .

Proof. As with classification, we divide the analysis into cases: the case in which d or n is constant, and the case in which each is sufficiently large.

Case d or n is a constant By our assumption that $\varepsilon^{-1} \leq \kappa \min(n, d)$, ε is a constant, and sufficiently large. So we just need to show an $\Omega(\min(M, n+d))$ bound. By the premise of the theorem, $M \geq 2(n+d)$, so we need only show an $\Omega(n+d)$ bound.

An $\Omega(d)$ bound. We construct an $n \times d$ matrix A as follows. For $i > 2$, each row A_i is just the vector $e_1 = (1, 0, 0, \dots, 0)$.

Let $A_{1,1} \leftarrow 0$, and initially assign $\zeta \leftarrow 1/d$ to all remaining entries of A_1 . Choose a random integer $j^* \in [2, d]$, and assign $A_{1,j^*} \leftarrow \sqrt{1-(d-2)\zeta^2}$. Note that $\|A_1\| = 1$.

Let $A_2 \leftarrow -A_1$, and then with probability $1/2$, negate A_{2,j^*} .

Our matrix A is as follows.

$$\begin{pmatrix} 0 & \zeta & \cdots & \zeta & \sqrt{1-(d-2)\zeta^2} & \zeta & \cdots & \zeta \\ 0 & -\zeta & \cdots & -\zeta & \pm\sqrt{1-(d-2)\zeta^2} & -\zeta & \cdots & -\zeta \\ 1 & 0 & \cdots & & & & & 0 \\ 1 & 0 & \cdots & & & & & 0 \\ 1 & \cdots & \cdots & & & & & \cdots \\ 1 & 0 & \cdots & & & & & 0 \end{pmatrix}$$

Observe that A has at most $2n + 2d \leq M$ non-zero entries, and all rows satisfy $\|A_i\| = 1$.

If $A_{1,j^*} = -A_{2,j^*}$, then A_1 and A_2 form a diametral pair, and the MEB radius is 1.

On the other hand, if $A_{1,j^*} = A_{2,j^*}$, then consider the ball center x with $x_1 \leftarrow x_{j^*} \leftarrow 1/\sqrt{2}$, and all other entries zero. Then for all $i > 2$, $\|x - A_i\|^2 = \left(1 - \frac{1}{\sqrt{2}}\right)^2$. On the other hand, for $i \in \{1, 2\}$, we have

$$\|x - A_i\|^2 \leq \frac{1}{2} + (d-2)\zeta^2 + \left(1 - \frac{1}{\sqrt{2}}\right)^2 \leq 2 - \sqrt{2} + \frac{1}{d}.$$

It follows that for ε satisfying the premise of the theorem, an algorithm satisfying the premise of the theorem can distinguish the two cases. By Fact 7.2, any algorithm distinguishing these two cases with probability $\geq 2/3$ must read $\Omega(d)$ entries of A .

An $\Omega(n)$ bound. We construct the $n \times d$ matrix A as follows. Initially set all rows $A_i \leftarrow e_1 = (1, 0, 0, \dots, 0)$. Then with probability $1/2$ choose a random $i^* \in [n]$, and negate $A_{i^*,1}$.

We have the following picture.

$$\begin{pmatrix} 1 & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 1 & 0 & \cdots & 0 \\ \pm 1 & 0 & \cdots & 0 \\ 1 & 0 & \cdots & 0 \\ \cdots & 0 & \cdots & 0 \\ 1 & 0 & \cdots & 0 \end{pmatrix}$$

The number of non-zeros of the resulting matrix is $n < M$. In the case where there is an entry of -1 , the MEB radius of A is 1, but otherwise the MEB radius is 0. Hence, an algorithm satisfying the premise of the theorem can distinguish the two cases. By Fact 7.2, any algorithm distinguishing these two cases with probability $\geq 2/3$ must read $\Omega(n)$ entries of A .

Case: d and n are sufficiently large Suppose first that $M = \Omega(\varepsilon^{-1}(n+d))$ for a sufficiently large constant in the $\Omega(\cdot)$. Put $s = \Theta(\varepsilon^{-1})$. We can assume without loss of generality that n , d , and s are sufficiently large integers. We need the following simple claim.

Claim 7.7. *Given an instance of the minimum enclosing ball problem in $T > t$ dimensions on a matrix with rows $\{\alpha e_i + \beta \sum_{j \in [t] \setminus \{i\}} e_j\}_{i=1}^t$ for distinct standard unit vectors e_i and $\alpha \geq \beta \geq 0$, the solution $x = \sum_{i=1}^t (\alpha + (t-1)\beta)e_i/t$ of cost $(\alpha - \beta)^2(1 - 1/t)$ is optimal.*

Proof. We can subtract the point $\beta \mathbf{1}_T$ from each of the points, and an optimal solution y for the translated problem yields an optimal solution $y + \beta \mathbf{1}_T$ for the original problem with the same cost. We can assume without loss of generality that $T = t$ and that e_1, \dots, e_t are the t standard unit

vectors in \mathbb{R}^t . Indeed, the value of each of the rows on each of the remaining coordinates is 0. The cost of the point $y_* = \sum_{i=1}^t (\alpha - \beta) e_i / t$ in the translated problem is

$$(\alpha - \beta)^2 \left(1 - \frac{1}{t}\right)^2 + (t - 1) (\alpha - \beta)^2 / t^2 = (\alpha - \beta)^2 \left(1 - \frac{1}{t}\right).$$

On the other hand, for any point y , the cost with respect to row i is $(\alpha - \beta - y_i)^2 + \sum_{j \neq i} (\beta - y_j)^2$. By averaging and Cauchy-Schwarz, there is a row of cost at least

$$\begin{aligned} \frac{1}{t} \cdot \left[\sum_{i=1}^t (\alpha - \beta - y_i)^2 + (t - 1) \sum_{i=1}^t y_i^2 \right] &= \|y\|^2 + (\alpha - \beta)^2 - \frac{2(\alpha - \beta) \|y\|_1}{t} \\ &\geq \|y\|^2 + (\alpha - \beta)^2 - \frac{2(\alpha - \beta) \|y\|}{\sqrt{t}} \end{aligned}$$

Taking the derivative w.r.t. to $\|y\|$, this is minimized when $\|y\| = \frac{\alpha - \beta}{\sqrt{t}}$, for which the cost is at least $(\alpha - \beta)^2 (1 - 1/t)$. \square

An $\Omega(ns)$ bound. We set the first s rows of A to e_1, \dots, e_s . We set all entries outside of the first s columns of A to 0. We choose the remaining $n - s = \Omega(n)$ rows of A by applying Theorem 7.1 with parameters $p = n - s, q = s$, and $W = 1/d$. If A is a YES instance, then by Claim 7.7, there is a solution with cost $(a - b)^2 (1 - 1/s) = 1 - \Theta(1/s)$. On the other hand, if A is a NO instance, then for a given x , either $\|A_{j^*} - x\|^2$ or $\|A_{p/2+j^*} - x\|^2$ is at least $a^2 = 1 - O(1/d)$. By setting $s = \Theta(\varepsilon^{-1})$ appropriately, these two cases differ by an additive ε , as needed.

An $\Omega(ds)$ bound. We choose A by applying Theorem 7.1 with parameters $p = s, q = d$, and $W = 1/d$. If A is a YES instance, then by Claim 7.7, there is a solution of cost at most $(a - b)^2 (1 - 1/s) = 1 - \Theta(1/s)$. On the other hand, if A is a NO instance, then for a given x , either $\|A_{j^*} - x\|^2$ or $\|A_{p/2+j^*} - x\|^2$ is at least $a^2 = 1 - O(1/d)$. As before, setting $s = \Theta(\varepsilon^{-1})$ appropriately causes these cases to differ by an additive ε .

Finally, it remains to show an $\Omega(M)$ bound in case $M = O(\varepsilon^{-1}(n + d))$. We will use our previous construction for showing an $\Omega(ns)$ bound, but replace the value of n there with n' , where n' is the largest integer for which $n's \leq M/2$. We claim that $n' \geq 1$. To see this, by the premise of the theorem $M \geq 2(n + d)$. Moreover, $s = \Theta(\varepsilon^{-1})$ and $\varepsilon^{-1} \leq \kappa(n + d)$. For a small enough constant $\kappa > 0$, $s \leq (n + d) \leq M/2$, as needed.

As the theorem statement concerns matrices with n rows, each of unit norm, we must have an input A with n rows. In this case, since the first row of A is e_1 , which has sparsity 1, we can simply set all remaining rows to the value of e_1 , without changing the MEB solution. Hence, the lower bound is $\Omega(n's) = \Omega(M)$. Notice that the number of non-zero entries is at most $n + n's \leq M/2 + M/2 = M$, as needed.

This completes the proof. \square

7.3 An $\tilde{\Omega}(n\varepsilon^{-2})$ Bound for Minimum Enclosing Ball

7.3.1 Intuition

Before diving into the intricate lower bound of this section, we describe a simple construction which lies at its core. Consider two distributions over arrays of size d : the first distribution, μ , is uniformly distributed over all strings with exactly $\frac{3d}{4}$ entries that are 1, and $\frac{d}{4}$ entries that are -1 . The second

distribution σ , is uniformly distributed over all strings with exactly $\frac{3d}{4} - D$ entries that are 1, and $\frac{d}{4} + D$ entries that are -1 , for $D = \tilde{O}(\sqrt{d})$.

Let $x \sim \mu$ with probability $\frac{1}{2}$ and $x \sim \sigma$ with probability $\frac{1}{2}$. Consider the task of deciding from which distribution x was sampled. In both cases, the distributions are over the sphere of radius \sqrt{d} , so the norm itself cannot be used to distinguish them. At the heart of our construction lies the following fact:

Fact 7.8. *Any algorithm that decides with probability $\geq \frac{3}{4}$ the distribution that x was sampled from, must read at least $\tilde{\Theta}(d)$ entries from x .*

We prove a version of this fact in the next sections. But first, let us explain the use of this fact in the lower bound construction: We create an instance of MEB which contains either n vectors similar to the first type, or alternatively $n - 1$ vector of the first type and an extra vector of the second type (with a small bias). To distinguish between the two types of instances, an algorithm has no choice but to check all n vectors, and for each invest $O(d)$ work as per the above fact. In our parameter setting, we'll choose $d = \tilde{O}(\varepsilon^{-2})$, attaining the lower bound of $\tilde{O}(nd) = \tilde{O}(n\varepsilon^{-2})$ in terms of time complexity.

To compute the difference in MEB center as $n \mapsto \infty$, note that by symmetry in the first case the center will be of the form (a, a, \dots, a) , where the value $a \in \mathbb{R}$ is chosen to minimize the maximal distance:

$$\arg \min_a \left\{ \frac{3}{4}(1-a)^2 + \frac{1}{4}(-1-a)^2 \right\} = \arg \min_a \{a^2 - a + 1\} = \frac{1}{2}$$

The second MEB center will be

$$\arg \min_a \left\{ \left(\frac{3}{4} - \frac{D}{d}\right)(1-a)^2 + \left(\frac{1}{4} + \frac{D}{d}\right)(-1-a)^2 \right\} = \arg \min_a \{a^2 - (1 - \frac{4D}{d})a + 1\} = \frac{1}{2} - \frac{2D}{d}$$

Hence, the difference in MEB centers is on the order of $\sqrt{d} \times (\frac{D}{d})^2 = O(D^2/d) = O(1)$. However, the whole construction is scaled to fit in the unit ball, and hence the difference in MEB centers becomes $\frac{1}{\sqrt{d}} \sim \varepsilon$. Hence for an ε approximation the algorithm must distinguish between the two distributions, which in turn requires $\Omega(\varepsilon^{-2})$ work.

7.3.2 Probabilistic Lemmas

For a set S of points in \mathbb{R}^d , let $\text{MEB}(S)$ denote the smallest ball that contains S . Let $\text{Radius}(S)$ be the radius of $\text{MEB}(S)$, and $\text{Center}(S)$ the unique center of $\text{MEB}(S)$.

For our next lower bound, our bad instance will come from points on the hypercube $\mathcal{H}_d = \{-\frac{1}{\sqrt{d}}, \frac{1}{\sqrt{d}}\}^d$.

Call a vertex of \mathcal{H}_d *regular* if it has $\frac{3d}{4}$ coordinates equal to $\frac{1}{\sqrt{d}}$ and $\frac{d}{4}$ coordinates equal to $-\frac{1}{\sqrt{d}}$. Call a vertex *special* if it has $\frac{3d}{4} - 12dD$ coordinates equal to $\frac{1}{\sqrt{d}}$ and $\frac{d}{4} + 12dD$ coordinates equal to $-\frac{1}{\sqrt{d}}$, where $D \equiv \frac{\ln n}{\sqrt{d}}$.

We will consider instances where all but one of the input rows A_i are random regular points, and one row may or may not be a random special point. We will need some lemmas about these points.

Lemma 7.9. *Let a denote a random regular point, b a special point, and c denote the point*

$\mathbf{1}_d/2\sqrt{d} = (\frac{1}{2\sqrt{d}}, \frac{1}{2\sqrt{d}}, \dots, \frac{1}{2\sqrt{d}})$. Then

$$\|a\|^2 = \|b\|^2 = 1 \quad (27)$$

$$\|c\|^2 = a^\top c = \frac{1}{4} \quad (28)$$

$$\|a - c\|^2 = \frac{3}{4} \quad (29)$$

$$b^\top c = \mathbf{E}[a^\top b] = \frac{1}{4} - 12D \quad (30)$$

Proof. The norm claims are entirely straightforward, and we have

$$a^\top c = \frac{1}{2d} \cdot \frac{3d}{4} - \frac{1}{2d} \cdot \frac{d}{4} = \frac{1}{4}.$$

Also (29) follows by

$$\|a - c\|^2 = \|a\|^2 + \|c\|^2 - 2a^\top c = 1 + \frac{1}{4} - 2\frac{1}{4} = \frac{3}{4}.$$

For (30), we have

$$b^\top c = \frac{1}{2d} \left(\frac{3d}{4} - 12dD \right) - \frac{1}{2d} \left(\frac{d}{4} + 12dD \right) = \frac{3}{8} - 6D - \frac{1}{8} - 6D = \frac{1}{4} - 12D,$$

and by linearity of expectation,

$$\begin{aligned} \mathbf{E}[a^\top b] &= d \cdot \frac{1}{d} \cdot \left(\frac{3}{4} \cdot \left(\frac{3}{4} - 12D \right) + \frac{1}{4} \cdot \left(\frac{1}{4} + 12D \right) - \frac{3}{4} \cdot \left(\frac{1}{4} + 12D \right) - \frac{1}{4} \cdot \left(\frac{3}{4} - 12D \right) \right) \\ &= \frac{1}{4} - 12D. \end{aligned}$$

□

Next, we show that $a^\top b$ is concentrated around its expectation (30).

Lemma 7.10. *Let a be a random regular point, and b a special point. For $d \geq 8 \ln^2 n$, $\Pr[a^\top b > \frac{1}{4} - 6D] \leq \frac{1}{n^3}$, and $\Pr[a^\top b < \frac{1}{4} - 18D] \leq \frac{1}{n^3}$.*

Proof. We will prove the first tail estimate, and then discuss the changes needed to prove the second estimate.

We apply the upper tail of the following enhanced form of Hoeffding's bound, which holds for random variables with bounded correlation.

Fact 7.11. *(Theorem 3.4 of [PS97] with their value of λ equal to 1) Let X_1, \dots, X_d be given random variables with support $\{0, 1\}$ and let $X = \sum_{j=1}^d X_j$. Let $\gamma > 0$ be arbitrary. If there exist independent random variables $\hat{X}_1, \dots, \hat{X}_d$ with $\hat{X} = \sum_{j=1}^d \hat{X}_j$ and $\mathbf{E}[X] \leq \mathbf{E}[\hat{X}]$ such that for all $J \subseteq [d]$,*

$$\Pr[\wedge_{j \in J} X_j = 1] \leq \prod_{j \in J} \Pr[\hat{X}_j = 1],$$

then

$$\Pr[X > (1 + \gamma) \mathbf{E}[\hat{X}]] \leq \left[\frac{e^\gamma}{(1 + \gamma)^{1 + \gamma}} \right]^{\mathbf{E}[\hat{X}]}.$$

Define $X_j = \frac{d}{2} \cdot (a_j b_j + \frac{1}{d})$. Since $a_j b_j \in \{-\frac{1}{\sqrt{d}}, \frac{1}{\sqrt{d}}\}$, the X_j have support $\{0, 1\}$. Let $\hat{X}_1, \dots, \hat{X}_d$ be i.i.d. variables with support $\{0, 1\}$ with $\mathbf{E}[\hat{X}_j] = \mathbf{E}[X_j]$ for all j .

We claim that for all $J \subseteq [d]$, $\Pr[\wedge_{j \in J} X_j = 1] \leq \prod_{j \in J} \Pr[\hat{X}_j = 1]$. By symmetry, it suffices to prove it for $J \in \{[1], [2], \dots, [d]\}$. We prove it by induction. The base case $J = [1]$ follows since $\mathbf{E}[\hat{X}_j] = \mathbf{E}[X_j]$. To prove the inequality for $J = [\ell]$, $\ell \geq 2$, assume the inequality holds for $[\ell - 1]$. Then,

$$\Pr[\wedge_{j \in [\ell]} X_j = 1] = \Pr[\wedge_{j \in [\ell-1]} X_j = 1] \cdot \Pr[X_\ell = 1 \mid \wedge_{j \in [\ell-1]} X_j = 1],$$

and by the inductive hypothesis,

$$\Pr[\wedge_{j \in [\ell-1]} X_j = 1] \leq \prod_{j \in [\ell-1]} \Pr[\hat{X}_j = 1],$$

so to complete the induction it is enough to show

$$\Pr[X_\ell = 1 \mid \wedge_{j \in [\ell-1]} X_j = 1] \leq \Pr[X_\ell = 1]. \quad (31)$$

Letting $\Delta(a, b)$ be the number of coordinates j for which $a_j \neq b_j$, we have

$$\Pr[X_\ell = 1] = 1 - \frac{\mathbf{E}[\Delta(a, b)]}{d}.$$

If $\wedge_{j \in [\ell-1]} X_j = 1$ occurs, then the first $\ell - 1$ coordinates of a_j and b_j have the same sign, and so

$$\Pr[X_\ell = 1 \mid \wedge_{j \in [\ell-1]} X_j = 1] = 1 - \frac{\mathbf{E}[\Delta(a, b) \mid \wedge_{j \in [\ell-1]} X_j = 1]}{d - \ell + 1} = 1 - \frac{\mathbf{E}[\Delta(a, b)]}{d - \ell + 1},$$

which proves (31).

We will apply Fact 7.11 to bound $\Pr[a^\top b > r]$ for $r = \frac{1}{4} - 6D$. Since $X = \frac{d}{2} a^\top b + \frac{d}{2} = \frac{d}{2}(1 + a^\top b)$, we have

$$\frac{X - \mathbf{E}[X]}{\mathbf{E}[X]} = \frac{a^\top b - \mathbf{E}[a^\top b]}{1 + \mathbf{E}[a^\top b]},$$

where we have used that (30) implies $\mathbf{E}[X]$ is positive (for large enough d), so we can perform the division. So

$$\frac{X - \mathbf{E}[X]}{\mathbf{E}[X]} - \frac{r - \mathbf{E}[a^\top b]}{1 + \mathbf{E}[a^\top b]} = \frac{a^\top b - \mathbf{E}[a^\top b]}{1 + \mathbf{E}[a^\top b]} - \frac{r - \mathbf{E}[a^\top b]}{1 + \mathbf{E}[a^\top b]} = \frac{a^\top b - r}{1 + \mathbf{E}[a^\top b]},$$

and so

$$\Pr[a^\top b > r] = \Pr\left[\frac{X - \mathbf{E}[X]}{\mathbf{E}[X]} > \frac{r - \mathbf{E}[a^\top b]}{1 + \mathbf{E}[a^\top b]}\right].$$

By Fact 7.11, for $\gamma = \frac{r - \mathbf{E}[a^\top b]}{1 + \mathbf{E}[a^\top b]}$, we have for $\gamma > 0$,

$$\Pr[a^\top b > r] \leq \left[\frac{e^\gamma}{(1 + \gamma)^{1 + \gamma}}\right]^{d(1 + \mathbf{E}[a^\top b])/2}.$$

By (30), $r - \mathbf{E}[a^\top b] = 6D$, and $1 \leq 1 + \mathbf{E}[a^\top b] \leq 2$, so $\gamma \in [3D, 6D]$. It is well-known (see Theorem 4.3 of [MR95]) that for $0 < \gamma < 2e - 1$, $e^\gamma \leq (1 + \gamma)^{1 + \gamma} e^{-\gamma^2/4}$, and so

$$\left[\frac{e^\gamma}{(1 + \gamma)^{1 + \gamma}}\right]^{d(1 + \mathbf{E}[a^\top b])/2} \leq \exp\left(-\frac{\gamma^2}{4}(d(1 + \mathbf{E}[a^\top b])/2)\right) = \exp(-\gamma^2 d(1 + \mathbf{E}[a^\top b])/8).$$

Since $\gamma \geq 3D$ and $\mathbf{E}[a^\top b] > 0$, this is at most $\exp(-D^2d) \leq \exp(-(\ln n)^2) \leq n^{-3}$, for large enough n , using the definition of D .

For the second tail estimate, we can apply the same argument to $-a$ and b , proving that $\Pr[-a^\top b > r] \leq 1/n^3$, where $r \equiv -1/4 + 18D$. We let X_j be the $\{0, 1\}$ variables $\frac{d}{2}(-a_j b_j + \frac{1}{d})$, with expected sum $\mathbf{E}[X] = 3d/8 + 6D$. As above, $\Pr[-a^\top b > r] = \Pr[\frac{X - \mathbf{E}[X]}{\mathbf{E}[X]} > \gamma]$, where $\gamma \equiv \frac{r - \mathbf{E}[-a^\top b]}{1 + \mathbf{E}[-a^\top b]}$. Now $\gamma\sqrt{d}$ is between $6 \ln n$ and $8 \ln n$, so the same relations apply as above, and the second tail estimate follows. \square

Note that since by (29) all regular points are distance $\sqrt{3}/2$ from c , that distance is an upper bound for the the MEB radius of a collection of regular points.

The next lemmas give more properties of MEBs involving regular and special points, under the assumption that the above concentration bounds on $a^\top b$ hold for a given special point b and all a in a collection of regular points.

That is, let S be a collection of random regular points. Let \mathcal{E} be the event that for all $a \in S$, $-18D \leq a^\top b - \frac{1}{4} \leq -6D$. By Lemma 7.10 and a union bound,

$$\Pr[\mathcal{E}] \geq 1 - \frac{2}{n^2},$$

when S has at most n points.

The condition of event \mathcal{E} applies not only to every point in S , but to every point in the convex hull $\text{conv } S$.

Lemma 7.12. *For special point b and collection S of points a , if event \mathcal{E} holds, then for every $a_S \in \text{conv } S$, $-18D \leq a_S^\top b - \frac{1}{4} \leq -6D$.*

Proof. Since $a_S \in \text{conv } S$, we have $a_S = \sum_{a \in S} p_a a$ for some values p_a with $\sum_{a \in S} p_a = 1$ and $p_a \geq 0$ for all $a \in S$. Therefore, assuming \mathcal{E} holds,

$$a_S^\top b = \left[\sum_{a \in S} p_a a \right]^\top b = \sum_{a \in S} p_a a^\top b \leq \sum_{a \in S} p_a (1/4 - 6D) = 1/4 - 6D,$$

and similarly $a_S^\top b \geq 1/4 - 18D$. \square

Lemma 7.13. *Suppose b is a special point and S is a collection of regular points such that event \mathcal{E} holds. Then for any $a_S \in \text{conv } S$, $\|a_S - b\| \geq \frac{\sqrt{3}}{2} + 6D$. Since $\text{Center}(S) \in \text{conv } S$, this bound applies to $\|\text{Center}(S) - b\|$ as well.*

Proof. Let H be the hyperplane normal to $c = \mathbf{1}_d/2\sqrt{d}$ and containing c . Then $S \subset H$, and so $\text{conv } S \subset H$, and since the minimum norm point in H is c , all points $a_S \in \text{conv } S$ have $\|a_S\|^2 \geq \|c\|^2 = 1/4$. By the assumption that event \mathcal{E} holds, and the previous lemma, we have $a_S^\top b \leq \frac{1}{4} - 6D$. Using this fact, $\|b\| = 1$, and $\|a_S\|^2 \geq 1/4$, we have

$$\begin{aligned} \|a_S - b\|^2 &= \|a_S\|^2 + \|b\|^2 - 2a_S^\top b \\ &\geq \frac{1}{4} + 1 - 2\left(\frac{1}{4} - 6D\right) \\ &= \frac{3}{4} + 12D, \end{aligned}$$

and so $\|a_S - b\| \geq \frac{\sqrt{3}}{2} + 6D$ provided D is smaller than a small constant. \square

Lemma 7.14. *Suppose a is a regular point, b is a special point, and $a^\top b \geq \frac{1}{4} - 18D$. Then there is a point $q \in \mathbb{R}^d$ for which $\|q - b\| = \frac{\sqrt{3}}{2}$ and $\|q - a\| \leq \frac{\sqrt{3}}{2} + \Theta(D^2)$, as $D \rightarrow 0$.*

Proof. As usual let $c \equiv \mathbf{1}_d/2\sqrt{d}$ and consider the point q at distance $\frac{\sqrt{3}}{2}$ from b on the line segment \overline{cb} , so

$$q = c + \gamma \cdot \frac{b - c}{\|b - c\|} = c + \gamma\alpha(b - c),$$

where $\alpha \equiv 1/\|b - c\|$ and γ is a value in $\Theta(D)$. From the definition of q ,

$$\begin{aligned} \|q - a\|^2 &= \|q\|^2 + \|a\|^2 - 2a^\top q \\ &= \|c\|^2 + 2\gamma\alpha b^\top c - 2\gamma\alpha\|c\|^2 + \gamma^2 + \|a\|^2 - 2a^\top c - 2\gamma\alpha a^\top b + 2\gamma\alpha a^\top c. \end{aligned}$$

Recall from (27) that $\|a\| = 1$, from (28) that $a^\top c = \|c\|^2 = \frac{1}{4}$, from (30) that $b^\top c = 1/4 - 12D$, and the assumption $a^\top b \geq 1/4 - 18D$, we have

$$\begin{aligned} \|q - a\|^2 &= 1/4 + 2\gamma\alpha(1/4 - 12D) - 2\gamma\alpha(1/4) + \gamma^2 + 1 - 2(1/4) - 2\gamma\alpha(1/4 - 18D) + 2\gamma\alpha(1/4) \\ &= 3/4 + 12\gamma\alpha D + \gamma^2 \\ &\leq 3/4 + \Theta(D^2), \end{aligned}$$

where the last inequality uses $\gamma = \Theta(D)$ and $\alpha = \Theta(1)$. □

7.3.3 Main Theorem

Given an $n \times d$ matrix A together with the norms $\|A_i\|$ for all rows A_i , as well as the promise that all $\|A_i\| = O(1)$, the ε -**MEB-Coreset** problem is to output a subset S of $\tilde{O}(\varepsilon^{-1})$ rows of A for which $A_i \in (1 + \varepsilon) \cdot \text{MEB}(S)$. Our main theorem in this section is the following.

Theorem 7.15. *If $n\varepsilon^{-1} \geq d$ and $d = \tilde{\Omega}(\varepsilon^{-2})$, then any randomized algorithm which with probability $\geq 4/5$ solves ε -**MEB-Coreset** must read $\tilde{\Omega}(n\varepsilon^{-2})$ entries of A for some choice of its random coins.*

We also define the following problem. Given an $n \times d$ matrix A together with the norms $\|A_i\|$ for all rows A_i , as well as the promise that all $\|A_i\| = O(1)$, the ε -**MEB-Center** problem is to output a vector $x \in \mathbb{R}^d$ for which $\|A_i - x\| \leq (1 + \varepsilon) \min_{y \in \mathbb{R}^d} \max_{i \in [n]} \|y - A_i\|$. We also show the following.

Theorem 7.16. *If $n\varepsilon^{-2} \geq d$ and $d = \tilde{\Omega}(\varepsilon^{-2})$, then any randomized algorithm which with probability $\geq 4/5$ solves ε -**MEB-Center** by outputting a convex combination of the rows A_i must read $\tilde{\Omega}(n\varepsilon^{-2})$ entries of A for some choice of its random coins.*

These theorems will follow from the same hardness construction, which we now describe. Put $d = 8\varepsilon^{-2} \ln^2 n$, which we assume is a sufficiently large power of 2. We also assume n is even. We construct two families \mathcal{F} and \mathcal{G} of $n \times d$ matrices A .

The family \mathcal{F} consists of all A for which each of the n rows in A is a regular point.

The family \mathcal{G} consists of all A for which exactly $n - 1$ rows of A are regular points, and one row of A is a special point.

(Recall that we say that a vertex of on \mathcal{H}_d is *regular* if it has exactly $\frac{3d}{4}$ coordinates equal to $\frac{1}{\sqrt{d}}$. We say a point on \mathcal{H}_d is *special* if it has exactly $d(\frac{3}{4} - 12D)$ coordinates equal to $\frac{1}{\sqrt{d}}$, where D is $\frac{\ln n}{\sqrt{d}}$.)

Let μ be the distribution on $n \times d$ matrices for which half of its mass is uniformly distributed on matrices in \mathcal{F} , while the remaining half is uniformly distributed on the matrices in \mathcal{G} . Let $\mathbf{A} \sim \mu$. We show that any randomized algorithm Alg which decides whether $\mathbf{A} \in \mathcal{F}$ or $\mathbf{A} \in \mathcal{G}$ with probability at least $3/4$ must read $\tilde{\Omega}(nd)$ entries of \mathbf{A} for some choice of its random coins. W.l.o.g., we may assume that Alg is deterministic, since we may average out its random coins, as we may fix its coin tosses that lead to the largest success probability (over the choice of \mathbf{A}). By symmetry and independence of the rows, we can assume that in each row, Alg queries entries in order, that is, if Alg makes s queries to a row A_i , we can assume it queries $A_{i,1}, A_{i,2}, \dots, A_{i,s}$, and in that order.

Let $r = d/(C \ln^2 n)$ for a sufficiently large constant $C > 0$. For a vector $u \in \mathbb{R}^d$, let $\text{pref}(u)$ denote its first r coordinates. Let ρ be the distribution of $\text{pref}(u)$ for a random regular point u . Let ρ' be the distribution of $\text{pref}(u)$ for a random special point u .

Lemma 7.17. (*Statistical Difference Lemma*) For $C > 0$ a sufficiently large constant,

$$\|\rho - \rho'\|_1 \leq \frac{1}{10}.$$

Proof. We will apply the following fact twice, once to ρ and once to ρ' .

Fact 7.18. (*special case of Theorem 4 of [DF80]*) Suppose an urn U contains d balls, each marked by one of two colors. Let H_{U^r} be the distribution of r draws made at random without replacement from U , and M_{U^r} be the distribution of r draws made at random with replacement. Then,

$$\|H_{U^r} - M_{U^r}\|_1 \leq \frac{4r}{d}.$$

Let σ be the distribution with support $\{\frac{1}{\sqrt{d}}, -\frac{1}{\sqrt{d}}\}$ with $\sigma(\frac{1}{\sqrt{d}}) = \frac{3}{4}$ and $\sigma(-\frac{1}{\sqrt{d}}) = \frac{1}{4}$. Let τ be the distribution with support $\{\frac{1}{\sqrt{d}}, -\frac{1}{\sqrt{d}}\}$ with $\tau(\frac{1}{\sqrt{d}}) = \frac{3}{4} - 12D$ and $\tau(-\frac{1}{\sqrt{d}}) = \frac{1}{4} + 12D$.

Let σ^r be the joint distribution of r independent samples from σ , and similarly define τ^r . Applying Fact 7.18 with $r = 1/100D^2$,

$$\|\rho - \sigma^r\|_1 \leq \frac{1}{25dD^2},$$

and

$$\|\rho' - \tau^r\|_1 \leq \frac{1}{25dD^2}.$$

By the triangle inequality,

$$\|\rho - \rho'\|_1 \leq \|\rho - \sigma^r\|_1 + \|\sigma^r - \tau^r\|_1 + \|\tau^r - \rho'\|_1 \leq \|\sigma^r - \tau^r\|_1 + \frac{2}{25dD^2},$$

and so it remains to bound $\|\sigma^r - \tau^r\|_1$. To do this, we use Stein's Lemma (see, e.g., ??, Section 12.8), which shows that for two coins with bias in $[\Omega(1), 1 - \Omega(1)]$, one needs $\Theta(z^{-2})$ independent coins tosses to distinguish the distributions with constant probability, where z is the difference in their expectations. Here, $z = 12D$, and so for constant $C > 0$ sufficiently large, for $r = 1/CD^2$, it follows that $\|\sigma^r - \tau^r\|_1 \leq \frac{1}{20}$. We thus have

$$\|\rho - \rho'\|_1 \leq \frac{1}{20} + \frac{2}{25dD^2} \leq \frac{1}{10},$$

where the last inequality uses $dD^2 = (\ln n)^2 \rightarrow \infty$. □

We use Lemma 7.17 to prove the following. We assume that Alg outputs 1 if it decides that $\mathbf{A} \in \mathcal{F}$, otherwise it outputs 0.

Theorem 7.19. *If Alg queries $o(nr)$ entries of \mathbf{A} , it cannot decide if $\mathbf{A} \in \mathcal{F}$ with probability at least $3/4$.*

Proof. We can think of \mathbf{A} as being generated according to the following random process.

1. Choose an index $i^* \in [n]$ uniformly at random.
2. Choose rows \mathbf{A}_j for $j \in [n]$ to be random independent regular points.
3. With probability $1/2$, do nothing. Otherwise, with the remaining probability $1/2$, replace \mathbf{A}_{i^*} with a random special point.
4. Output \mathbf{A} .

Define the advantage $\text{adv}(Alg)$ to be:

$$\text{adv}(Alg) \equiv \left| \Pr_{A \in_R \mathcal{G}} [Alg(A) = 1] - \Pr_{A \in_R \mathcal{F}} [Alg(A) = 1] \right|.$$

To prove the theorem, it suffices to show $\text{adv}(Alg) < 1/4$. Let $\bar{\mathbf{A}}_{i^*}$ denote the rows of \mathbf{A} , excluding row i^* , generated in step 2. By the description of the random process above, we have

$$\text{adv}(Alg) = \mathbf{E}_{i^*, \bar{\mathbf{A}}_{i^*}} \left| \Pr_{\text{special } \mathbf{A}_{i^*}} [Alg(A) = 1 \mid i^*, \bar{\mathbf{A}}_{i^*}] - \Pr_{\text{regular } \mathbf{A}_{i^*}} [Alg(A) = 1 \mid i^*, \bar{\mathbf{A}}_{i^*}] \right|.$$

To analyze this quantity, we first condition on a certain event $\mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*})$ holding, which will occur with probability $1 - o(1)$, and allow us to discard the pairs $(i^*, \bar{\mathbf{A}}_{i^*})$ that do not satisfy the condition of the event. Intuitively, the event is just that for most regular \mathbf{A}_{i^*} , algorithm Alg does not read more than r entries in \mathbf{A}_{i^*} . This holds with probability $1 - o(1)$, over the choice of i^* and $\bar{\mathbf{A}}_{i^*}$, because all n rows of A are i.i.d., and so on average Alg can only afford to read $o(r)$ entries in each row.

More formally, we say a pair $(i^*, \bar{\mathbf{A}}_{i^*})$ is *good* if

$$\Pr_{\text{regular } \mathbf{A}_{i^*}} [Alg \text{ queries at most } r \text{ queries of } \mathbf{A}_{i^*} \mid (i^*, \bar{\mathbf{A}}_{i^*}) = (i^*, \bar{\mathbf{A}}_{i^*})].$$

Let $\mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*})$ be the event that $(i^*, \bar{\mathbf{A}}_{i^*})$ is good. Then, $\Pr_{i^*, \bar{\mathbf{A}}_{i^*}} [\mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*})] = 1 - o(1)$, and we can upper bound the advantage by

$$\mathbf{E}_{i^*, \bar{\mathbf{A}}_{i^*}} \left| \Pr_{\text{special } \mathbf{A}_{i^*}} [Alg(A) = 1 \mid \mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*}), i^*, \bar{\mathbf{A}}_{i^*}] - \Pr_{\text{regular } \mathbf{A}_{i^*}} [Alg(A) = 1 \mid \mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*}), i^*, \bar{\mathbf{A}}_{i^*}] \right| + o(1).$$

Consider the algorithm Alg'_{i^*} , which on input A , makes the same sequence of queries to A as Alg unless it must query more than r positions of \mathbf{A}_{i^*} . In this case, it outputs an arbitrary value in $\{0, 1\}$, otherwise it outputs $Alg(A)$.

Claim 7.20.

$$\left| \Pr_{\text{regular } \mathbf{A}_{i^*}} [Alg(A) = 1 \mid \mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*}), i^*, \bar{\mathbf{A}}_{i^*}] - \Pr_{\text{regular } \mathbf{A}_{i^*}} [Alg'_{i^*}(A) = 1 \mid i^*, \bar{\mathbf{A}}_{i^*}] \right| = o(1),$$

Proof. Since $\mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*})$ occurs,

$$\Pr_{\text{regular } \mathbf{A}_{i^*}} [\text{Alg makes at most } r \text{ queries to } \mathbf{A}_{i^*} \mid i^*, \bar{\mathbf{A}}_{i^*}] = 1 - o(1).$$

This implies that

$$\left| \Pr_{\text{regular } \mathbf{A}_{i^*}} [\text{Alg}(A) = 1 \mid \mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*}), i^*, \bar{\mathbf{A}}_{i^*}] - \Pr_{\text{regular } \mathbf{A}_{i^*}} [\text{Alg}'_{i^*}(A) = 1 \mid i^*, \bar{\mathbf{A}}_{i^*}] \right| = o(1).$$

□

By Lemma 7.17, we have that

$$\left| \Pr_{\text{regular } \mathbf{A}_{i^*}} [\text{Alg}'_{i^*}(A) = 1 \mid i^*, \bar{\mathbf{A}}_{i^*}] - \Pr_{\text{special } \mathbf{A}_{i^*}} [\text{Alg}'_{i^*}(A) = 1 \mid i^*, \bar{\mathbf{A}}_{i^*}] \right| \leq \frac{1}{10}.$$

Hence, by Claim 7.20 and the triangle inequality, we have that

$$\left| \Pr_{\text{regular } \mathbf{A}_{i^*}} [\text{Alg}(A) = 1 \mid \mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*}), i^*, \bar{\mathbf{A}}_{i^*}] - \Pr_{\text{special } \mathbf{A}_{i^*}} [\text{Alg}'_{i^*}(A) = 1 \mid i^*, \bar{\mathbf{A}}_{i^*}] \right| \leq \frac{1}{10} + o(1).$$

To finish the proof, it suffices to show the following claim

Claim 7.21.

$$\left| \Pr_{\text{special } \mathbf{A}_{i^*}} [\text{Alg}(A) = 1 \mid \mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*}), i^*, \bar{\mathbf{A}}_{i^*}] - \Pr_{\text{special } \mathbf{A}_{i^*}} [\text{Alg}'_{i^*}(A) = 1 \mid i^*, \bar{\mathbf{A}}_{i^*}] \right| \leq \frac{1}{10} + o(1).$$

Indeed, if we show Claim 7.21, then by the triangle inequality we will have that $\text{adv}(\text{Alg}) \leq \frac{1}{5} + o(1) < \frac{1}{4}$.

Proof of Claim 7.21: Since $\mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*})$ occurs,

$$\Pr_{\text{regular } \mathbf{A}_{i^*}} [\text{Alg makes at most } r \text{ queries to } \mathbf{A}_{i^*} \mid i^*, \bar{\mathbf{A}}_{i^*}] = 1 - o(1).$$

Since ρ is the distribution of prefixes of regular points, this condition can be rewritten as

$$\Pr_{u \sim \rho} [\text{Alg makes at most } r \text{ queries to the } i^*\text{-th row} \mid i^*, \bar{\mathbf{A}}_{i^*}, \text{pref}(A_{i^*}) = u] = 1 - o(1).$$

By Lemma 7.17, we thus have,

$$\Pr_{u \sim \rho'} [\text{Alg makes at most } r \text{ queries to the } i^*\text{-th row} \mid i^*, \bar{\mathbf{A}}_{i^*}, \text{pref}(A_{i^*}) = u] \geq \frac{9}{10} - o(1).$$

Since ρ' is the distribution of prefixes of special points, this condition can be rewritten as

$$\Pr_{\text{special } \mathbf{A}_{i^*}} [\text{Alg makes at most } r \text{ queries to } \mathbf{A}_{i^*} \mid i^*, \bar{\mathbf{A}}_{i^*}] \geq \frac{9}{10} - o(1).$$

This implies that

$$\left| \Pr_{\text{special } \mathbf{A}_{i^*}} [\text{Alg}(A) = 1 \mid \mathcal{E}(i^*, \bar{\mathbf{A}}_{i^*}), i^*, \bar{\mathbf{A}}_{i^*}] - \Pr_{\text{special } \mathbf{A}_{i^*}} [\text{Alg}'_{i^*}(A) = 1 \mid i^*, \bar{\mathbf{A}}_{i^*}] \right| \leq \frac{1}{10} + o(1).$$

□

This completes the proof of the theorem. □

7.3.4 Proofs of Theorem 7.15 and 7.16

Next we show how Theorem 7.19 implies Theorem 7.15 and Theorem 7.16, using the results on MEBs of regular and special points.

Proof of Theorem 7.15: We set the dimension $d = 4 \cdot 36 \cdot \varepsilon^{-2} \ln^2(n-1)$. Let A' denote the set of regular rows of \mathbf{A} . We condition on event \mathcal{E} , namely, that every convex combination $p^T A$, where $p \in \Delta_{n-1}$, satisfies $p^T A' b \leq \frac{1}{4} - 6D$. This event occurs with probability at least $1 - 2n^{-2}$. (We may neglect the difference between n and $n-1$ in some expressions.)

It follows by Lemma 7.13 that if $\mathbf{A} \in \mathcal{G}$, then for every $S \subseteq A'$,

$$\|\text{Center}(S) - b\| \geq \frac{\sqrt{3}}{2} + 2\varepsilon.$$

By (29), $\text{Radius}(A') \leq \frac{\sqrt{3}}{2}$. It follows that any algorithm that, with probability at least $4/5$, outputs a subset S of $\tilde{O}(\varepsilon^{-1})$ rows of \mathbf{A} for which $\mathbf{A}_i \in (1 + \varepsilon) \cdot \text{MEB}(S)$ must include the point $b \in S$.

Given such an algorithm, by reading each of the $\tilde{O}(\varepsilon^{-1})$ rows output, we can determine if $\mathbf{A} \in \mathcal{F}$ or $\mathbf{A} \in \mathcal{G}$ with an additional $\tilde{O}(\varepsilon^{-1}d)$ time. By Theorem 7.19, the total time must be $\tilde{\Omega}(n\varepsilon^{-2})$. By assumption, $n\varepsilon^{-1} \geq d$, and so any randomized algorithm that solves ε -**MEB-Coreset** with probability at least $4/5$, can decide if $\mathbf{A} \in \mathcal{F}$ with probability at least $4/5 - 2n^{-2} \geq 3/4$, and so it must read $\tilde{\Omega}(n\varepsilon^{-2})$ entries for some choice of its random coins. \square

Proof of Theorem 7.16: We again set the dimension $d = 4 \cdot 36 \cdot \varepsilon^{-2} \ln^2(n-1)$. Let A' denote the set of regular rows of \mathbf{A} . We again condition on the event \mathcal{E} .

By Lemma 7.13, if $\mathbf{A} \in \mathcal{G}$, then for every convex combination $p^T A'$,

$$\|p^T A' - b\| \geq \frac{\sqrt{3}}{2} + 2\varepsilon,$$

and so the MEB radius returned by any algorithm that outputs a convex combination of rows of A' must be at least $\frac{\sqrt{3}}{2} + 2\varepsilon$.

However, by (29), if $\mathbf{A} \in \mathcal{F}$, then $\text{Radius}(\mathbf{A}) \leq \frac{\sqrt{3}}{2}$. On the other hand, by Lemma 7.14, if $\mathbf{A} \in \mathcal{G}$, then $\text{MEB-radius}(\mathbf{A}) \leq \frac{\sqrt{3}}{2} + \Theta(\varepsilon^2)$.

It follows that if $\mathbf{A} \in \mathcal{G}$, then the convex combination $p^T \mathbf{A}$ output by the algorithm must have a non-zero coefficient multiplying the special point b . This, in particular, implies that $p^T \mathbf{A}$ is not on the affine hyperplane H with normal vector $\mathbf{1}_d$ containing the point $c = \mathbf{1}_d/2\sqrt{d}$. However, if $\mathbf{A} \in \mathcal{F}$, then any convex combination of the points is on H . The output $p^T \mathbf{A}$ of the algorithm is on H if and only if $p^T \mathbf{A} \mathbf{1}_d = \frac{\sqrt{d}}{2}$, which can be tested in $O(d)$ time.

By Theorem 7.19, the total time must be $\tilde{\Omega}(n\varepsilon^{-2})$. By assumption, $n\varepsilon^{-2} \geq d$, and so any randomized algorithm that solves ε -**MEB-Center** with probability $\geq 4/5$ by outputting a convex combination of rows can decide if $\mathbf{A} \in \mathcal{F}$ with probability at least $4/5 - 2n^{-2} \geq 3/4$, and so must read $\tilde{\Omega}(n\varepsilon^{-2})$ entries for some choice of its random coins. \square

7.4 Las Vegas Algorithms

While our algorithms are Monte Carlo, meaning they err with small probability, it may be desirable to obtain Las Vegas algorithms, i.e., randomized algorithms that have low expected time but never err. We show this cannot be done in sublinear time.

Theorem 7.22. *For the classification and minimum enclosing ball problems, there is no Las Vegas algorithm that reads an expected $o(M)$ entries of its input matrix and solves the problem to within a one-sided additive error of at most $1/2$. This holds even if $\|A_i\| = 1$ for all rows A_i .*

Proof. Suppose first that $n \geq M$. Consider $n \times d$ matrices A, B^1, \dots, B^M , where for each $C \in \{A, B^1, \dots, B^M\}$, $C_{i,j} = 0$ if either $j > 1$ or $i > M$. Also, $A_{i,1} = 1$ for $i \in [M]$, while for each j , $B_{1,i}^j = 1$ if $i \in [M] \setminus \{j\}$, while $B_{1,j}^j = -1$. With probability $1/2$ the matrix A is chosen, otherwise a matrix B^j is chosen for a random j . Notice that whichever case we are in, each of the first M rows of the input matrix has norm equal to 1, while all remaining rows have norm 0. It is easy to see that distinguishing these two cases with probability $\geq 2/3$ requires reading $\Omega(M)$ entries. As $\Omega(M)$ is a lower bound for Monte Carlo algorithms, it is also a lower bound for Las Vegas algorithms. Moreover, distinguishing these two cases is necessary, since if the problem is classification, if $C = A$ the margin is 1, otherwise it is 0, while if the problem is minimum enclosing ball, if $C = A$ the cost is 0, otherwise it is 1.

We now assume $M > n$. Let d' be the largest integer for which $nd' < M$. Here $d' \geq 1$. Let A be the $n \times d'$ matrix, where $A_{i,j} = \frac{1}{\sqrt{d'}}$ for all i and j . The margin of A is 1, and the minimum enclosing ball has radius 0.

Suppose there were an algorithm Alg on input A for which there is an assignment to Alg 's random tape r for which Alg reads at most $nd'/4$ of its entries. If there were no such r , the expected running time of Alg is already $\Omega(nd') = \Omega(M)$. Let A_ℓ be a row of A for which Alg reads at most $d'/4$ entries of A_ℓ given random tape r , and let $S \subset [d']$ be the set of indices in A_ℓ read, where $|S| \leq d'/4$. Consider the $n \times d'$ matrix B for which $B_{i,j} = A_{i,j}$ for all $i \neq \ell$, while $B_{\ell,j} = A_{\ell,j}$ for all $j \in S$, and $B_{\ell,j} = -A_{\ell,j}$ for all $j \in [d'] \setminus S$. Notice that all rows of A and B have norm 1.

To bound the margin of B , consider any vector x of norm at most 1. Then

$$(A_\ell + B_\ell)x \leq \|x\| \cdot \|A_\ell + B_\ell\| \leq \|A_\ell + B_\ell\|.$$

$A_\ell + B_\ell$ has at least $3d'/4$ entries that are 0, while the non-zero entries all have value $2/\sqrt{d'}$. Hence, $\|A_\ell + B_\ell\|^2 \leq \frac{d'}{4} \cdot \frac{4}{d'} = 1$. It follows that either $A_\ell x$ or $B_\ell x$ is at most $1/2$, which bounds the margin of B . As Alg cannot distinguish A and B given random tape r , it cannot have one-sided additive error at most $1/2$.

For minimum enclosing ball, notice that $\|A_\ell - B_\ell\|^2 \cdot \frac{1}{4} \geq \frac{3d'}{4} \cdot \frac{4}{d'} \cdot \frac{1}{4} = \frac{3}{4}$, which lower bounds the cost of the minimum enclosing ball of B . As Alg cannot distinguish A and B given random tape r , it cannot have one-sided additive error at most $3/4$. \square

8 Concluding Remarks

We have described a general method for sublinear optimization of constrained convex programs, and showed applications to classical problems in machine learning such as linear classification and minimum enclosing ball obtaining improvements in leading-order terms over the state of the art. The application of our sublinear primal-dual algorithms to soft margin SVM and related convex problems is currently explored in ongoing work with Nati Srebro.

In all our running times the dimension d can be replaced by the parameter S , which is the maximum over the input rows A_i of the number of nonzero entries in A_i . Note that $d \geq S \geq M/n$. Here we require the assumption that entries of any given row can be recovered in $O(S)$ time, which is compatible with keeping each row as a hash table or (up to a logarithmic factor in run-time) in sorted order.

Acknowledgements We thank Nati Srebro and an anonymous referee for helpful comments on the relation between this work and PAC learning theory.

References

- [AS10] Pankaj Agarwal and R. Sharathkumar. Streaming algorithms for extent problems in high dimensions. In *SODA '10: Proc. Twenty-First ACM-SIAM Symposium on Discrete Algorithms*, 2010.
- [BFKV98] Avrim Blum, Alan M. Frieze, Ravi Kannan, and Santosh Vempala. A polynomial-time algorithm for learning noisy linear threshold functions. *Algorithmica*, 22(1/2):35–52, 1998.
- [Byl94] Tom Bylander. Learning linear threshold functions in the presence of classification noise. In *COLT '94: Proceedings of the Seventh Annual Conference on Computational Learning Theory*, pages 340–347, New York, NY, USA, 1994. ACM.
- [CBCG04] Nicolò Cesa-Bianchi, Alex Conconi, and Claudio Gentile. On the generalization ability of on-line learning algorithms. *IEEE Transactions on Information Theory*, 50(9):2050–2057, 2004.
- [Cla08] Kenneth L. Clarkson. Coresets, sparse greedy approximation, and the Frank-Wolfe algorithm. In *SODA '08: Proc. Nineteenth ACM-SIAM Symposium on Discrete Algorithms*, pages 922–931, Philadelphia, PA, USA, 2008. Society for Industrial and Applied Mathematics.
- [DF80] P. Diaconis and D. Freedman. Finite exchangeable sequences. *The Annals of Probability*, 8:745–764, 1980.
- [DV04] John Dunagan and Santosh Vempala. A simple polynomial-time rescaling algorithm for solving linear programs. In *STOC '04: Proceedings of the Thirty-Sixth Annual ACM Symposium on the Theory of Computing*, pages 315–320, New York, NY, USA, 2004. ACM.
- [FKM⁺08] Joan Feigenbaum, Sampath Kannan, Andrew McGregor, Siddharth Suri, and Jian Zhang. Graph distances in the data-stream model. *SIAM J. Comput.*, 38(5):1709–1727, 2008.
- [FW56] Marguerite Frank and Philip Wolfe. An algorithm for quadratic programming. *Naval Res. Logist. Quart.*, 3:95–110, 1956.
- [GK95] Michael D. Grigoriadis and Leonid G. Khachiyan. A sublinear-time randomized approximation algorithm for matrix games. *Operations Research Letters*, 18:53–58, 1995.
- [Haz10] E. Hazan. The convex optimization approach to regret minimization. In *to appear in Optimization for Machine Learning*, available at <http://www.cs.princeton.edu/~ehazan/papers/OCO-survey.pdf>. MIT Press, 2010.
- [HKKA06] Elad Hazan, Adam Kalai, Satyen Kale, and Amit Agarwal. Logarithmic regret algorithms for online convex optimization. In Gábor Lugosi and Hans-Ulrich Simon, editors, *COLT*, volume 4005 of *Lecture Notes in Computer Science*, pages 499–513. Springer, 2006.

- [KY07] Christos Koufogiannakis and Neal E. Young. Beating simplex for fractional packing and covering linear programs. In *FOCS*, pages 494–504. IEEE Computer Society, 2007.
- [MP88] M. L. Minsky and S. Papert. *Perceptrons: An introduction to computational geometry*. MIT press Cambridge, Mass, 1988.
- [MR95] Rajeev Motwani and Prabakar Raghavan. *Randomized Algorithms*. Cambridge University Press, 1995.
- [Mut05] S. Muthukrishnan. Data streams: Algorithms and applications. *Foundations and Trends in Theoretical Computer Science*, 1(2), 2005.
- [MW10] Morteza Monemizadeh and David Woodruff. 1-pass relative error l_p -sampling with applications. In *SODA '10: Proc. Twenty-First ACM-SIAM Symposium on Discrete Algorithms*, 2010.
- [Nov62] A.B.J. Novikoff. On convergence proofs on perceptrons. In *Proceedings of the Symposium on the Mathematical Theory of Automata, Vol XII*, pages 615–622, 1962.
- [PS97] Alessandro Panconesi and Aravind Srinivasan. Randomized distributed edge coloring via an extension of the chernoff-hoeffding bounds. *SIAM J. Comput.*, 26(2):350–368, 1997.
- [PST91] Serge A. Plotkin, David B. Shmoys, and Éva Tardos. Fast approximation algorithms for fractional packing and covering problems. In *SFCS '91: Proceedings of the 32nd Annual Symposium on Foundations of Computer Science*, pages 495–504, Washington, DC, USA, 1991. IEEE Computer Society.
- [Ser99] Rocco A. Servedio. On PAC learning using winnow, perceptron, and a perceptron-like algorithm. In *COLT '99: Proceedings of the Twelfth Annual Conference on Computational Learning Theory*, pages 296–307, New York, NY, USA, 1999. ACM.
- [SS03] Bernhard Schölkopf and Alexander J. Smola. A short introduction to learning with kernels. pages 41–64, 2003.
- [SV09] Ankan Saha and S.V.N. Vishwanathan. Efficient approximation algorithms for minimum enclosing convex shapes. arXiv:0909.1062v2, 2009.
- [TZ04] Mikkel Thorup and Yin Zhang. Tabulation based 4-universal hashing with applications to second moment estimation. In *SODA '04: Proc. Fifteenth ACM-SIAM Symposium on Discrete Algorithms*, pages 615–624, Philadelphia, PA, USA, 2004. Society for Industrial and Applied Mathematics.
- [Zin03] Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the Twentieth International Conference on Machine Learning(ICML)*, pages 928–936, 2003.
- [ZZC06] Hamid Zarrabi-Zadeh and Timothy M. Chan. A simple streaming algorithm for minimum enclosing balls. In *CCCG*, 2006.

A Main Tools

A.1 Tools from online learning

Online linear optimization The following lemma is essentially due to Zinkevich [Zin03]:

Lemma A.1 (OGD). *Consider a set of vectors $q_1, \dots, q_T \in \mathbb{R}^d$ such that $\|q_i\|_2 \leq c$. Let $x_0 \leftarrow 0$, and $\tilde{x}_{t+1} \leftarrow x_t + \frac{1}{\sqrt{T}}q_t$, $x_{t+1} \leftarrow \frac{\tilde{x}_{t+1}}{\max\{1, \|\tilde{x}_{t+1}\|\}}$. Then*

$$\max_{x \in \mathbb{B}} \sum_{t=1}^T q_t^\top x - \sum_{t=1}^T q_t^\top x_t \leq 2c\sqrt{T}.$$

This is true even if each q_t is dependent on x_1, \dots, x_{t-1} .

Proof. Assume $c = 1$, generalization is by straightforward scaling. Let $\eta = \frac{1}{\sqrt{T}}$. By definition and for any $\|x\| \leq 1$,

$$\|x - x_{t+1}\|^2 \leq \|x - \tilde{x}_{t+1}\|^2 = \|x - x_t - \eta q_t\|^2 = \|x - x_t\|^2 - 2\eta q_t^\top (x - x_t) + \eta^2 \|q_t\|^2.$$

Rearranging we obtain

$$q_t^\top (x - x_t) \leq \frac{1}{2\eta} [\|x - x_t\|^2 - \|x - x_{t+1}\|^2] + \eta/2.$$

Summing up over $t = 1$ to T yields

$$\sum_t q_t^\top x - \sum_t q_t^\top x_t \leq \frac{1}{2\eta} \|x - x_1\|^2 + \eta T/2 \leq \frac{2}{\eta} + \frac{\eta}{2} T \leq 2\sqrt{T}.$$

□

For our streaming and parallel implementation, a simpler version of gradient descent, also essentially due to Zinkevich [Zin03], is given by:

Lemma A.2 (Lazy Projection OGD). *Consider a set of vectors $q_1, \dots, q_T \in \mathbb{R}^d$ such that $\|q_i\|_2 \leq 1$. Let*

$$x_{t+1} \leftarrow \arg \min_{x \in \mathbb{B}} \left\{ \sum_{\tau=1}^t q_\tau^\top \cdot x + \sqrt{2T} \|x\|_2^2 \right\}$$

Then

$$\max_{x \in \mathbb{B}} \sum_{t=1}^T q_t^\top x - \sum_{t=1}^T q_t^\top x_t \leq 2\sqrt{2T}.$$

This is true even if each q_t is dependent on x_1, \dots, x_{t-1} .

For a proof see Theorem 2.1 in [Haz10], where we take $\mathcal{R}(x) = \|x\|_2^2$, and the norm of the linear cost functions is bounded by $\|q_t\|_2 \leq 1$, as is the diameter of \mathcal{K} - the ball in our case. Notice that the solution of the above optimization problem is simply:

$$x_{t+1} = \frac{y_{t+1}}{\max\{1, \|y_{t+1}\|\}}, \quad y_{t+1} = -\frac{\sum_{\tau=1}^t q_\tau}{\sqrt{2T}}$$

Strongly convex loss functions The following Lemma is essentially due to [HKKA06].

For $H \in \mathbb{R}$ with $H > 0$, a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is H -strongly convex in \mathbb{B} if for all $x \in \mathbb{B}$, all the eigenvalues of $\nabla^2 f(x)$ are at least H .

Lemma A.3 (OGDStrictlyConvex). *Consider a set of H -strongly convex functions f_1, \dots, f_T such that the norm of their gradients is bounded over the unit ball \mathbb{B} by $G \geq \max_t \max_{x \in \mathbb{B}} \|\nabla f_t(x)\|$. Let $x_0 \in \mathbb{B}$, and $\tilde{x}_{t+1} \leftarrow x_t - \frac{1}{t} \nabla f_t(x_t)$, $x_{t+1} \leftarrow \frac{\tilde{x}_{t+1}}{\max\{1, \|\tilde{x}_{t+1}\|\}}$. Then*

$$\sum_{t=1}^T f_t(x_t) - \min_{\|x\|_2 \leq 1} \sum_{t=1}^T f_t(x) \leq \frac{G^2}{H} \log T.$$

This is true even if each f_t is dependent on x_1, \dots, x_{t-1} .

Again, for the MEB application and its relatives it is easier to implement the lazy versions in the streaming model. The following Lemma is the analogous tool we need:

Lemma A.4. *Consider a set of H -strongly convex functions f_1, \dots, f_T such that the norm of their gradients is bounded over the unit ball \mathbb{B} by $G \geq \max_t \max_{x \in \mathbb{B}} \|\nabla f_t(x)\|$. Let*

$$x_{t+1} \leftarrow \arg \min_{x \in \mathbb{B}} \left\{ \sum_{\tau=1}^t f_\tau(x) \right\}$$

Then

$$\sum_{t=1}^T f_t(x_t) - \min_{x \in \mathbb{B}} \sum_{t=1}^T f_t(x) \leq \frac{2G^2}{H} \log T.$$

This is true even if each f_t is dependent on x_1, \dots, x_{t-1} .

Proof. By Lemma 2.3 in [Haz10] we have:

$$\sum_{t=1}^T f_t(x_t) - \min_{\|x\|_2 \leq 1} \sum_{t=1}^T f_t(x) \leq \sum_t [f_t(x_t) - f_t(x_{t+1})]$$

Denote by $\Phi_t(x) = \sum_{\tau=1}^t f_\tau$. Then by Taylor expansion at x_{t+1} , there exists a $z_t \in [x_{t+1}, x_t]$ for which

$$\begin{aligned} \Phi_t(x_t) &= \Phi_t(x_{t+1}) + (x_t - x_{t+1})^\top \nabla \Phi_t(x_{t+1}) + \frac{1}{2} \|x_t - x_{t+1}\|_{z_t}^2 \\ &\geq \Phi_t(x_{t+1}) + \frac{1}{2} \|x_t - x_{t+1}\|_{z_t}^2, \end{aligned}$$

using the notation $\|y\|_z^2 = y^\top \nabla^2 \Phi_t(z) y$. The inequality above is true because x_{t+1} is a minimum of Φ_t over \mathcal{K} . Thus,

$$\begin{aligned} \|x_t - x_{t+1}\|_{z_t}^2 &\leq 2 \Phi_t(x_t) - 2 \Phi_t(x_{t+1}) \\ &= 2 (\Phi_{t-1}(x_t) - \Phi_{t-1}(x_{t+1})) + 2[f_t(x_t) - f_t(x_{t+1})] \\ &\leq 2[f_t(x_t) - f_t(x_{t+1})] && \text{optimality of } x_t \\ &\leq 2 \nabla f_t(x_t)^\top (x_t - x_{t+1}) && \text{convexity of } f_t. \end{aligned}$$

By convexity and Cauchy-Schwarz:

$$\begin{aligned} f_t(x_t) - f_t(x_{t+1}) &\leq \nabla f_t(x_t)(x_t - x_{t+1}) \leq \|\nabla f_t(x_t)\|_{z_t}^* \|x_t - x_{t+1}\|_{z_t} \\ &\leq \|\nabla f_t(x_t)\|_{z_t}^* \sqrt{2 \nabla f_t(x_t)^\top (x_t - x_{t+1})} \end{aligned}$$

Shifting sides and squaring, we get

$$f_t(x_t) - f_t(x_{t+1}) \leq \nabla f_t(x_t)(x_t - x_{t+1}) \leq 2\|\nabla f_t(x_t)\|_{z_t}^{*2}$$

Since f_t are assumed to be H -strongly convex, we have $\|\cdot\|_z \geq \|\cdot\|_{Ht}$, and hence for the dual norm,

$$f_t(x_t) - f_t(x_{t+1}) \leq 2\|\nabla f_t(x_t)\|_{z_t}^{*2} \leq 2 \frac{\|\nabla f_t(x_t)\|_2^2}{Ht} \leq \frac{2G^2}{Ht}$$

Summing over all iterations we get

$$\sum_{t=1}^T f_t(x_t) - \min_{\|x\|_2 \leq 1} \sum_{t=1}^T f_t(x) \leq \sum_t [f_t(x_t) - f_t(x_{t+1})] \leq \sum_t \frac{2G^2}{Ht} \leq \frac{2G^2}{H} \log T$$

□

Combining sampling and regret minimization

Lemma A.5. Consider a set of H -strongly convex functions f_1, \dots, f_T such that the norm of their gradients is bounded over the unit ball by $G \geq \max_t \max_{x \in \mathbb{B}} \|\nabla f_t(x)\|$. Let

$$y_{t+1} \leftarrow \begin{cases} \arg \min_{x \in \mathbb{B}} \{\sum_{\tau=1}^t f_\tau(x)\} & \text{w.p. } \alpha \\ y_t & \text{o/w} \end{cases}$$

Then for a fixed x^* we have

$$\mathbf{E}[\sum_{t=1}^T f_t(y_t) - \sum_{t=1}^T f_t(x^*)] \leq \frac{1}{\alpha} \frac{2G^2}{H} \log T.$$

This is true even if each f_t is dependent on y_1, \dots, y_{t-1} .

Proof. Consider the sequence of functions \tilde{f}_t defined as

$$\tilde{f}_t \leftarrow \begin{cases} \frac{f_t}{\alpha} & \text{w.p. } \alpha \\ 0 & \text{o/w} \end{cases}$$

Where 0 denotes the all-zero function. Then the algorithm from Lemma A.4 applied to the functions \tilde{f}_t is exactly the algorithm we apply above to the functions f_t . Notice that the functions \tilde{f}_t are $\frac{H}{\alpha}$ -strongly convex, and in addition their gradients are bounded by $\frac{G}{\alpha}$. Hence applying Lemma A.4 we obtain

$$\mathbf{E}[\sum_{t=1}^T f_t(y_t) - \sum_{t=1}^T f_t(x^*)] = \mathbf{E}[\sum_{t=1}^T \tilde{f}_t(x_t) - \sum_{t=1}^T \tilde{f}_t(x^*)] \leq \frac{1}{\alpha} \frac{2G^2}{H} \log T.$$

□

B Auxiliary lemmas

First, some simple lemmas about random variables.

Lemma B.1. *Let X be a random variable with $|\mathbf{E}[X]| \leq C$, and let $\bar{X} = \text{clip}(X, C) = \min\{C, \max\{-C, X\}\}$ for some $C \in \mathbb{R}$. Then*

$$|\mathbf{E}[\bar{X}] - \mathbf{E}[X]| \leq \frac{\mathbf{Var}[X]}{C}.$$

Proof. By direct calculation:

$$\begin{aligned} \mathbf{E}[\bar{X}] - \mathbf{E}[X] &= \int_{x < -C} \Pr[x](-C - x) + \int_{x > C} \Pr[x](C - x), \\ &\leq \int_{x < -C} \Pr[x]|x| - \int_{x < -C} \Pr[x]C \\ &\leq \int_{x < -C} \Pr[x]x^2/C - \int_{x < -C} \Pr[x]C \\ &= \int_{x < -C} \Pr[x] \frac{x^2 - C^2}{C} \\ &\leq \int_{x < -C} \Pr[x] \frac{x^2 - \mathbf{E}[X]^2}{C} && \text{since } |\mathbf{E}[X]| \leq C \\ &= \frac{\mathbf{Var}[X^2]}{C} \end{aligned}$$

and similarly $\mathbf{E}[\bar{X}] - \mathbf{E}[X] \geq -\mathbf{Var}[X]/C$, and the result follows. □

Lemma B.2. *For random variables X and Y , and $\alpha \in [0, 1]$,*

$$\mathbf{E}[(\alpha X + (1 - \alpha)Y)^2] \leq \max\{\mathbf{E}[X^2], \mathbf{E}[Y^2]\}.$$

This implies by induction that the second moment of a convex combination of random variables is no more than the maximum of their second moments.

Proof. We have, using Cauchy-Schwarz for the first inequality,

$$\begin{aligned} \mathbf{E}[(\alpha X + (1 - \alpha)Y)^2] &= \alpha^2 \mathbf{E}[X^2] + 2\alpha(1 - \alpha) \mathbf{E}[XY] + (1 - \alpha)^2 \mathbf{E}[Y^2] \\ &\leq \alpha^2 \mathbf{E}[X^2] + 2\alpha(1 - \alpha) \sqrt{\mathbf{E}[X^2] \mathbf{E}[Y^2]} + (1 - \alpha)^2 \mathbf{E}[Y^2] \\ &= (\alpha \sqrt{\mathbf{E}[X^2]} + (1 - \alpha) \sqrt{\mathbf{E}[Y^2]})^2 \\ &\leq \max\{\sqrt{\mathbf{E}[X^2]}, \sqrt{\mathbf{E}[Y^2]}\}^2 \\ &= \max\{\mathbf{E}[X^2], \mathbf{E}[Y^2]\}. \end{aligned}$$

□

B.1 Martingale and concentration lemmas

The Bernstein inequality, that holds for random variables $Z_t, t \in [T]$ that are independent, and such that for all t , $\mathbf{E}[Z_t] = 0$, $\mathbf{E}[Z_t^2] \leq s$, and $|Z_t| \leq V$, states

$$\log \text{Prob}\left\{ \sum_{t \in [T]} Z_t \geq \alpha \right\} \leq -\alpha^2/2(Ts + \alpha V/3) \tag{32}$$

Here we need a similar bound for random variables which are not independent, but form a martingale with respect to a certain filtration. Many concentration results have been proven for Martingales, including somewhere, in all likelihood, the present lemma. However, for clarity and completeness, we will outline how the proof of the Bernstein inequality can be adapted to this setting.

Lemma B.3. *Let $\{Z_t\}$ be a martingale difference sequence with respect to filtration $\{S_t\}$, such that $\mathbf{E}[Z_t|S_1, \dots, S_t] = 0$. Assume the filtration $\{S_t\}$ is such that the values in S_t are determined using only those in S_{t-1} , and not any previous history, and so the joint probability distribution*

$$\text{Prob}\{S_1 = s_1, S_2 = s_2, \dots, S_T = s_T\} = \prod_{t \in [T-1]} \text{Prob}\{S_{t+1} = s_{t+1} \mid S_t = s_t\},$$

In addition, assume for all t , $\mathbf{E}[Z_t^2|S_1, \dots, S_t] \leq s$, and $|Z_t| \leq V$. Then

$$\log \text{Prob}\left\{\sum_{t \in T} Z_t \geq \alpha\right\} \leq -\alpha^2/2(Ts + \alpha V/3).$$

Proof. A key step in proving the Bernstein inequality is to show an upper bound on the exponential generating function $\mathbf{E}[\exp(\lambda Z)]$, where $Z \equiv \sum_t Z_t$, and $\lambda > 0$ is a parameter to be chosen. This step is where the hypothesis of independence is applied. In our setting, we can show a similar upper bound on this expectation: Let $\mathbf{E}_t[\cdot]$ denote expectation with respect to S_t , and $\mathbf{E}_{[T]}$ denote expectation with respect to S_t for $t \in [T]$. This expression for the probability distribution implies that for any real-valued function f of state tuples S_t ,

$$\begin{aligned} \mathbf{E}_{[T]}[\prod_{t \in [T]} f(S_t)] &= f(s_1) \int_{s_2, \dots, s_T} [\prod_{t \in [T-1]} f(s_{t+1})][\prod_{t \in [T-1]} \text{Prob}\{S_{t+1} = s_{t+1} \mid S_t = s_t\}] \\ &= f(s_1) \int_{s_2, \dots, s_{T-1}} \left[[\prod_{t \in [T-2]} f(s_{t+1})][\prod_{t \in [T-2]} \text{Prob}\{S_{t+1} = s_{t+1} \mid S_t = s_t\}] \right. \\ &\quad \left. \int_{s_T} f(s_T) \text{Prob}\{S_T = s_T \mid S_{T-1} = s_{T-1}\} \right], \end{aligned}$$

where the inner integral can be denoted as the conditional expectation $\mathbf{E}_T[f(S_T) \mid S_{T-1}]$. By induction this is

$$f(s_1) \left[\int_{s_2} f(s_2) \text{Prob}\{S_2 = s_2 \mid S_1 = s_1\} \left[\int_{s_3} \dots \int_{s_T} f(s_T) \text{Prob}\{S_T = s_T \mid S_{T-1} = s_{T-1}\} \dots \right] \right],$$

and by writing the constant $f(S_1)$ as the expectation with respect to the constant $S_0 = s_0$, and using $\mathbf{E}_X[\mathbf{E}_X[Y]] = \mathbf{E}_X[Y]$ for any random variables X and Y , we can write this as

$$\mathbf{E}_{[T]}[\prod_{t \in [T]} f(S_t)] = \mathbf{E}_{[T]}[\prod_{t \in [T]} \mathbf{E}_t[f(S_t) \mid S_{t-1}]].$$

For fixed i and a given $\lambda \in \mathbb{R}$, we take $f(S_1) = 1$, and $f(S_t) \equiv \exp(\lambda Z_{t-1})$, to obtain

$$\mathbf{E}_{[T]} \left[\exp\left(\lambda \sum_{t \in [T]} Z_t\right) \right] = \mathbf{E}_{[T]} \left[\prod_{t \in [T]} \mathbf{E}_t[\exp(\lambda Z_t) \mid S_{t-1}] \right].$$

Now for *any* random variable X with $\mathbf{E}[X] = 0$, $\mathbf{E}[X^2] \leq s$, and $|X| \leq V$,

$$\mathbf{E}[\exp(\lambda X)] \leq \exp\left(\frac{s}{V^2}(e^{\lambda V} - 1 - \lambda V)\right),$$

(as is shown and used for proving Bernstein's inequality in the independent case) and therefore

$$\mathbf{E}_{[T]}[\exp(\lambda Z)] \leq \mathbf{E}_{[T]}\left[\prod_{t \in [T]} \exp\left(\frac{s}{V^2}(e^{\lambda V} - 1 - \lambda V)\right)\right] = \exp\left(T \frac{s}{V^2}(e^{\lambda V} - 1 - \lambda V)\right).$$

where $Z \equiv \sum_{t \in [T]} Z_t$. This bound is the same as is obtained for independent Z_t , and so the remainder of the proof is exactly as in the proof for the independent case: Markov's inequality is applied to the random variable $\exp(\lambda Z)$, obtaining

$$\text{Prob}\{Z \geq \alpha\} \leq \exp(-\lambda \alpha) \mathbf{E}_{[T]}[\exp(\lambda Z)] \leq \exp(-\lambda \alpha + T \frac{s}{V^2}(e^{\lambda V} - 1 - \lambda V)),$$

and an appropriate value $\lambda = \frac{1}{V} \log(1 + \alpha V/Ts)$ is chosen for minimizing the bound, yielding

$$\text{Prob}\{Z \geq \alpha\} \leq \exp\left(-\frac{Ts}{V^2}((1 + \gamma) \log(1 + \gamma) - \gamma)\right),$$

where $\gamma \equiv \alpha V/Ts$, and finally the inequality for $\gamma \geq 0$ that $(1 + \gamma) \log(1 + \gamma) - \gamma \geq \frac{\gamma^2/2}{1 + \gamma/3}$ is applied. \square

B.2 Proof of lemmas used in main theorem

We restate and prove lemmas 2.4, 2.5 and 2.6, in slightly more general form. In the following we only assume that $v_t(i) = \text{clip}(\tilde{v}_t(i), \frac{1}{\eta})$ is the clipping of a random variable $\tilde{v}_t(i)$. The variance of $\tilde{v}_t(i)$ is at most one $\mathbf{Var}[\tilde{v}_t(i)] \leq 1$, and we denote by $\mu_t(i) = \mathbf{E}[\tilde{v}_t(i)]$. We also assume that the expectations of $\tilde{v}_t(i)$ are bounded by an absolute constant $|\mu_t(i)| \leq C \leq \frac{1}{\eta}$. This constant is one for the perceptron application, but at most two for MEB. Note that since the variance of $\tilde{v}_t(i)$ is bounded by one, so is the variance of its clipping $v_t(i)$ ².

Lemma B.4. For $\eta \leq \sqrt{\frac{\log n}{10T}}$, with probability at least $1 - O(1/n)$,

$$\max_i \sum_{t \in [T]} [v_t(i) - \mu_t(i)] \leq 90\eta T.$$

Proof. Lemma B.1 implies that $|\mathbf{E}[v_t(i)] - \mu_t(i)| \leq \eta$, since $\mathbf{Var}[\tilde{v}_t(i)] \leq 1$.

We show that for given $i \in [n]$, with probability $1 - O(1/n^2)$, $\sum_{t \in [T]} [v_t(i) - \mathbf{E}[v_t(i)]] \leq 80\eta T$, and then apply the union bound over all $i \in [n]$. This together with the above bound on $|\mathbf{E}[v_t(i)] - \mu_t(i)|$ implies the lemma via the triangle inequality.

Fixing i , let $Z_t^i \equiv v_t(i) - \mathbf{E}[v_t(i)]$, and consider the filtration given by

$$S_t \equiv (x_t, p_t, w_t, y_t, v_{t-1}, i_{t-1}, j_{t-1}, v_{t-1} - \mathbf{E}[v_{t-1}]),$$

Using the notation $\mathbf{E}_t[\cdot] = \mathbf{E}[\cdot | S_t]$, Observe that

²This follows from the fact that the second moment only decreases by the clipping operation, and definition of variance as $\mathbf{Var}(v_t(i)) = \min_z \mathbf{E}[v_t(i)^2 - z^2]$. We can use $z = \mathbf{E}[\tilde{v}_t(i)]$, and hence the decrease in second moment suffices.

1. $\forall t . \mathbf{E}_t[(Z_t^i)^2] = \mathbf{E}_t[v_t(i)^2] - \mathbf{E}_t[v_t(i)]^2 = \mathbf{Var}(v_t(i)) \leq 1.$
2. $|Z_t^i| \leq 2/\eta.$ This holds since by construction, $|v_t(i)| \leq 1/\eta,$ and hence

$$|Z_t^i| = |v_t(i) - \mathbf{E}[v_t(i)]| \leq |v_t(i)| + |\mathbf{E}[v_t(i)]| \leq \frac{2}{\eta}$$

Using these conditions, despite the fact that the Z_t^i are not independent, we can use Lemma B.3, and conclude that $Z \equiv \sum_{t \in T} Z_t^i$ satisfies the Bernstein-type inequality with $s = 1$ and $V = 2/\eta$

$$\log \text{Prob}\{Z \geq \alpha\} \leq -\alpha^2/2(Ts + \alpha V/3) \leq -\alpha^2/2(T + 2\alpha/3\eta),$$

Letting $\alpha \leftarrow 80\eta T,$ we have

$$\log \text{Prob}\{Z \geq 80\eta T\} \leq -\alpha^2/2(T + 2\alpha/3\eta) \leq -20\eta^2 T$$

For $\eta = \sqrt{\frac{\log n}{10T}},$ above probability is at most $e^{-2\log n} \leq \frac{1}{n^2}.$ □

Lemma 2.5 can be restated in the following more general form:

Lemma B.5. For $\eta \leq \sqrt{\frac{\log n}{10T}},$ with probability at least $1 - O(1/n),$ it holds that $\left| \sum_{t \in [T]} \mu_t(i_t) - \sum_t p_t^\top v_t \right| \leq 100C\eta T.$

It is a corollary of the following two lemmas:

Lemma B.6. For $\eta \leq \sqrt{\frac{\log n}{10T}},$ with probability at least $1 - O(1/n),$

$$\left| \sum_{t \in [T]} p_t^\top v_t - \sum_t p_t^\top \mu_t \right| \leq 90\eta T.$$

Proof. This Lemma is proven in essentially the same manner as Lemma 2.4, and proven below for completeness.

Lemma B.1 implies that $|\mathbf{E}[v_t(i)] - \mu_t(i)| \leq \eta,$ using $\mathbf{Var}[\tilde{v}_t(i)] \leq 1.$ Since p_t is a distribution, it follows that $|\mathbf{E}[p_t^\top v_t] - p_t^\top \mu_t| \leq \eta$

Let $Z_t \equiv p_t^\top v_t - \mathbf{E}[p_t^\top v_t] = \sum_i p_t(i) Z_t^i,$ where $Z_t^i = v_t(i) - \mathbf{E}[v_t(i)].$ Consider the filtration given by

$$S_t \equiv (x_t, p_t, w_t, y_t, v_{t-1}, i_{t-1}, j_{t-1}, v_{t-1} - \mathbf{E}[v_{t-1}]),$$

Using the notation $\mathbf{E}_t[\cdot] = \mathbf{E}[\cdot | S_t],$ the quantities $|Z_t|$ and $\mathbf{E}_t[Z_t^2]$ can be bounded as follows:

$$|Z_t| = \left| \sum_i p_t(i) Z_t^i \right| \leq \sum_i p_t(i) |Z_t^i| \leq 2\eta^{-1} \quad \text{using } |Z_t^i| \leq 2\eta^{-1} \text{ as in Lemma 2.4.}$$

Also, using properties of variance, we have

$$\mathbf{E}[Z_t^2] = \text{Var}[p_t^\top v_t] = \sum_i p_t(i)^2 \mathbf{Var}(v_t(i)) \leq \max_i \mathbf{Var}[v_t(i)] \leq 1.$$

We can now apply the Bernstein-type inequality of Lemma B.3, and continue exactly as in Lemma 2.4. □

Lemma B.7. For $\eta \leq \sqrt{\frac{\log n}{10T}}$, with probability at least $1 - O(1/n)$,

$$\left| \sum_{t \in [T]} \mu_t(i_t) - \sum_t p_t \mu_t \right| \leq 10C\eta T.$$

Proof. Let $Z_t \equiv \mu_t(i_t) - p_t \mu_t$, where now μ_t is a constant vector and i_t is the random variable, and consider the filtration given by

$$S_t \equiv (x_t, p_t, w_t, y_t, v_{t-1}, i_{t-1}, j_{t-1}, Z_{t-1}),$$

The expectation of $\mu_t(i_t)$, conditioning on S_t with respect to the random choice $r(i_t)$, is $p_t \mu_t$. Hence $\mathbf{E}_t[Z_t] = 0$, where $\mathbf{E}_t[\cdot]$ denotes $\mathbf{E}[\cdot | S_t]$. The parameters $|Z_t|$ and $\mathbf{E}[Z_t^2]$ can be bounded as follows:

$$|Z_t| \leq |\mu_t(i)| + |p_t \mu_t| \leq 2C$$

$$\mathbf{E}[Z_t^2] = \mathbf{E}[(\mu_t(i) - p_t^\top \mu_t)^2] \leq 2\mathbf{E}[\mu_t(i)^2] + 2(p_t^\top \mu_t)^2 \leq 4C^2$$

Applying Lemma B.3 to $Z \equiv \sum_{t \in T} Z_t$, with parameters $s \leq 4C^2$, $V \leq 2C$, we obtain

$$\log \text{Prob}\{Z \geq \alpha\} \leq -\alpha^2 / (4C^2 T + 2C\alpha),$$

Letting $\alpha \leftarrow 10C\eta T$, we obtain

$$\log \text{Prob}\{Z \geq 10\eta T\} \leq -\frac{100\eta^2 C^2 T^2}{4C^2 T + 20C^2 \eta T} \leq 5\eta^2 T \leq \log n$$

Where the last inequality holds assuming $\eta \leq \sqrt{\frac{\log n}{T}}$. □

Finally, we prove Lemma 2.6 by a simple application of Markov's inequality:

Lemma B.8. *w.p. at least $1 - \frac{1}{4}$ it holds that $\sum_t p_t^\top v_t^2 \leq 8C^2 T$.*

Proof. By assumption, $\mathbf{E}[v_t^2(i)] \leq C^2$, and using Lemma B.1, we have $\mathbf{E}[v_t(i)^2] \leq (C + \frac{1}{C})^2 \leq 2C^2$.

By linearity of expectation, we have $\mathbf{E}[\sum_t p_t^\top v_t^2] \leq 2C^2 T$, and since the random variables v_t^2 are non-negative, applying Markov's inequality yields the lemma. □

C Bounded precision

All algorithms in this paper can be implemented with bounded precision.

First we observe that approximation of both the training data and the vectors that are “played” does not increase the regret too much, for both settings we are working in.

Lemma C.1. *Given a sequence of functions f_1, \dots, f_T and another sequence $\tilde{f}_1, \dots, \tilde{f}_T$ all mapping \mathbb{R}^d to \mathbb{R} , such that $|\tilde{f}_t(x) - f_t(x)| \leq \alpha_f$ for all $x \in B$ and $t \in [T]$, suppose $x_1, \dots, x_T \in B$ is a sequence of regret R against $\{f_t\}$, that is,*

$$\max_{x \in B} \sum_{t \in [T]} \tilde{f}_t(x) - \sum_{t \in [T]} \tilde{f}_t(x_t) \leq R.$$

Now suppose $\tilde{x}_1, \dots, \tilde{x}_T \in \mathbb{R}^d$ is a sequence with $|f_t(\tilde{x}_t) - f_t(x_t)| \leq \alpha_x$ for all $t \in [T]$. Then

$$\max_{x \in B} \sum_{t \in [T]} f_t(x) - \sum_{t \in [T]} f_t(\tilde{x}_t) \leq R + T(\alpha_x + 2\alpha_f).$$

Proof. For $x \in \mathbb{B}$, we have $\sum_{t \in [T]} f_t(x) \leq \sum_{t \in [T]} \tilde{f}_t(x) + T\alpha_f$, and

$$\sum_{t \in [T]} f_t(\tilde{x}_t) \geq \sum_{t \in [T]} f_t(x_t) - T\alpha_x \geq \sum_{t \in [T]} \tilde{f}_t(x_t) - T\alpha_x - T\alpha_f,$$

and the result follows by combining these inequalities. \square

That is, x_t is some sequence known to have small regret against the “training functions” $\tilde{f}_t(x)$, which are approximations to the true functions of interest, and the \tilde{x}_t are approximations to these x_t . The lemma says that despite these approximations, the \tilde{x}_t sequence has controllable regret against the true functions.

This lemma is stated in more generality than we need: all functions considered here have the form $f_t(x) = b_t + q_t^\top x + \gamma \|x\|^2$, where $|b_t| \leq 1$, $q_t \in \mathbb{B}$, and $|\gamma| \leq 1$. Thus if $\tilde{f}_t(x) = \tilde{b}_t + \tilde{q}_t^\top x + \gamma \|x\|^2$, then the first condition $|\tilde{f}_t(x) - f_t(x)| \leq \alpha_f$ holds when $|b_t - \tilde{b}_t| + \|q_t - \tilde{q}_t\| \leq \alpha_f$. Also, the second condition $|f_t(\tilde{x}_t) - f_t(x_t)| \leq \alpha_x$ holds for such functions when $\|\tilde{x}_t - x_t\| \leq \alpha_x/3$.

Lemma C.2. *Given a sequence of vectors $q_1, \dots, q_T \in \mathbb{R}^n$, with $\|q_t\|_\infty \leq B$ for $t \in [T]$, and a sequence $\tilde{q}_1, \dots, \tilde{q}_T \in \mathbb{R}^n$ such that $\|\tilde{q}_t - q_t\|_\infty \leq \alpha_q$ for all $t \in [T]$, suppose $p_1, \dots, p_T \in \Delta$ is a sequence of regret R against $\{\tilde{q}_t\}$, that is,*

$$\sum_{t \in [T]} p_t^\top \tilde{q}_t - \min_{p \in \Delta} \sum_{t \in [T]} p^\top \tilde{q}_t \leq R.$$

Now suppose $\tilde{p}_1, \dots, \tilde{p}_T \in \mathbb{R}^n$ is a sequence with $\|\tilde{p}_t - p_t\|_1 \leq \alpha_p$ for all $t \in [T]$. Then

$$\sum_{t \in [T]} \tilde{p}_t^\top q_t - \min_{p \in \Delta} \sum_{t \in [T]} p^\top q_t \leq R + T(B\alpha_p + 2\alpha_q).$$

Proof. For $p \in \Delta$ we have $\sum_{t \in [T]} p^\top q_t \geq \sum_{t \in [T]} p^\top \tilde{q}_t + T\alpha_q$, and

$$\sum_{t \in [T]} \tilde{p}_t^\top q_t \leq \sum_{t \in [T]} p_t^\top q_t + TB\alpha_p \leq \sum_{t \in [T]} p_t^\top \tilde{q}_t + TB\alpha_p + T\alpha_q,$$

The proof follows by combining the inequalities. \square

Note that to have $\|\tilde{p}_t - p_t\|_1 \leq \alpha_p$, it is enough that the relative error of each entry of \tilde{p}_t is α_p .

The use of \tilde{q}_t in place of q_t (for either of the two lemmas) will be helpful for our semi-streaming and kernelized algorithms (§5, §6), where computation of the norms $\|y_t\|$ of the working vectors y_t is a bottleneck; the above two lemmas imply that it is enough to compute such norms to within relative ϵ or so.

C.1 Bit Precision for Algorithm 1

First, the bit precision needed for the OGD part of the algorithm. Let γ denote a sufficiently small constant fraction of ϵ , where the small constant is absolute. From Lemma C.1 and following discussion, we need only use the rows A_i up to a precision that gives an approximation \tilde{A}_i that is within Euclidean distance γ , and similarly for an approximation \tilde{x}_t of x_t . For the latter, in particular, we need only compute $\|y_t\|$ to within relative error γ . Thus a per-entry precision of γ/\sqrt{d} is sufficient.

We need $\|x_t\|$ for ℓ_2 sampling; arithmetic relative error γ/\sqrt{d} in the sampling procedure gives an estimate of $\tilde{v}_t(i)$ for which $\mathbf{E}[A\tilde{v}_t] = A\hat{x}_t$, where \hat{x}_t is a vector within $O(\gamma)$ Euclidean distance of x_t . We can thus charge this error to the OGD analysis, where \hat{x}_t is the \tilde{x}_t of Lemma C.1.

For the MW part of the algorithm, we observe that due to the clipping step, if the initial computation of $\tilde{v}_t(i)$, Line 9, is done with $\eta\epsilon/5$ relative error, then the computed value is within $\epsilon/5$ additive error. Similar precision for the clipping implies that the computed value of $v_t(i)$, which takes the place of \tilde{q}_t in Lemma C.2, is within $\epsilon/5$ of the exact version, corresponding to q_t in the lemma. Here B of the lemma, bounding $\|q_t\|_\infty$, is $1/\eta$, due to the clipping.

It remains to determine the arithmetic relative error needed in the update step, Line 11, to keep the relative error of the computed value of p_t , or \tilde{p}_t of Lemma C.2, small enough. Indeed, if the relative error is a small enough constant fraction of $\eta\epsilon/T$, then the relative error of all updates together can be $\eta\epsilon/3$. Thus $\alpha_p \leq \eta\epsilon/3$ and $\alpha_q \leq \epsilon/3$ and the added regret due to arithmetic error is at most $T\epsilon$.

Summing up: the arithmetic precision needed is at most on the order of

$$-\log \min\{\epsilon/\sqrt{d}, \eta\epsilon, \eta\epsilon/T\} = O(\log(nd/\epsilon)),$$

to obtain a solution with additive $T\epsilon/10$ regret over the solution computed using exact computation. This implies an additional error of $\epsilon/10$ to the computed solution, and thus changes only constant factors in the algorithm.

C.2 Bit Precision for Convex Quadratic Programming

From the remarks following Lemma C.1, the conditions of that lemma hold in the setting of convex quadratic programming in the simplex, assuming that every $A_i \in \mathbb{B}$. Thus the discussion of §C.1 carries over, up to constants, with the simplification that computation of $\|y_t\|$ is not needed.