

# Worst-case Quadratic Loss Bounds for Prediction Using Linear Functions and Gradient Descent

Nicolò Cesa-Bianchi, Philip M. Long, and Manfred K. Warmuth

*Abstract*— In this paper we study the performance of gradient descent when applied to the problem of on-line linear prediction in arbitrary inner product spaces. We prove worst-case bounds on the sum of the squared prediction errors under various assumptions concerning the amount of *a priori* information about the sequence to predict. The algorithms we use are variants and extensions of on-line gradient descent. Whereas our algorithms always predict using linear functions as hypotheses, none of our results requires the data to be linearly related. In fact, the bounds proved on the total prediction loss are typically expressed as a function of the total loss of the best fixed linear predictor with bounded norm. All the upper bounds are tight to within constants. Matching lower bounds are provided in some cases. Finally, we apply our results to the problem of on-line prediction for classes of smooth functions.

*Keywords*— prediction, Widrow-Hoff algorithm, gradient descent, adaptive linear filter theory, smoothing, inner product spaces, computational learning theory, on-line learning, linear systems, worst-case loss bounds.

## I. INTRODUCTION

In this paper we analyze algorithms in the on-line prediction model. This model was introduced by Angluin [1] and Littlestone [16], [17]. Unlike other settings, where the predictor’s goal is to estimate a set of parameters in a nearly optimal way with respect to some criterion, the goal in this model is to generate predictions, in a sequential fashion, so as to minimize the *total (sum) loss* over the whole sequence of examples. Throughout this paper we use the squared prediction error as the loss function for each example. This loss is sometimes called the square loss. Though we focus on the performance of on-line algorithms from a purely theoretical viewpoint, one of the main contributions of this study is the derivation of the optimal learning rate

Nicolò Cesa-Bianchi is with DSI, Università di Milano, Via Comelico 39, 20135 Milano, ITALY. Part of this research was done while this author was visiting UC Santa Cruz partially supported by the “Progetto finalizzato sistemi informatici e calcolo parallelo” of CNR under grant 91.00884.69.115.09672 (Italy). Email address: [cesabian@dsi.unimi.it](mailto:cesabian@dsi.unimi.it).

Phil Long is with the Department of Computer Science, Duke University, P.O. Box 90129, Durham, NC 27708 USA. He is supported by ONR grant N00014-94-1-0938. Part of this work was done while he was supported by AFOSR grant F49620-92-J-051, and part was done while he was at Technische Universitaet Graz supported by a Lise Meitner Fellowship from the Fonds zur Förderung der wissenschaftlichen Forschung (Austria), and at UC Santa Cruz supported by a UCSC Chancellor’s dissertation-year fellowship. Email address: [plong@csc.duke.edu](mailto:plong@csc.duke.edu).

Manfred Warmuth is with the Computer Science Department, UC Santa Cruz, Santa Cruz, CA 95064 USA, and is supported by ONR grant N00014-91-J-1162 and NSF grant IRI-9123692. Email: [manfred@cse.ucsc.edu](mailto:manfred@cse.ucsc.edu).

for gradient descent applied to linear predictors.

We assume the prediction process occurs in a sequence of trials. At trial number  $t$  the prediction algorithm

- is presented with an *instance*  $x_t$  chosen from some domain  $\mathcal{X}$ ,
- is required to return a real number  $\hat{y}_t$ ,
- then receives a real number  $y_t$  from the environment which we interpret as the truth.

The total loss of an algorithm over a sequence of  $m$  trials is  $\sum_{t=1}^m (\hat{y}_t - y_t)^2$ . A critical aspect of this model is that when the algorithm is making its prediction  $\hat{y}_t$  for the  $t$ th instance  $x_t$ , it has access to pairs  $(x_s, y_s)$  *only for*  $s < t$ .

We adopt a worst-case outlook, following [27], [21], [18], [7], [22], [3] and many others, assuming nothing about the environment of the predictor, in particular the pairs  $(x_1, y_1), \dots, (x_m, y_m)$ . Our results can be loosely interpreted as having the following message: “To the extent that the environment is friendly, our algorithms have small total loss.” Of course, the strength of such results depends on how “friendly” is formalized. For the most general results of this paper (described in Section IV), the domain  $\mathcal{X}$  is assumed to be a real vector space.<sup>1</sup> To formalize “friendly,” we make use of the general notion of an inner product  $(\cdot, \cdot)$ , which is any function from  $\mathcal{X} \times \mathcal{X}$  to  $\mathbf{R}$  that has certain properties (see Section III for a list). The inner product formalization is very general. One of the simplest inner products may be defined as follows in the case that  $\mathcal{X} = \mathbf{R}^n$ :

$$(\mathbf{u}, \mathbf{v}) = \sum_{i=1}^n u_i v_i = \mathbf{u} \cdot \mathbf{v}.$$

Notice that for any inner product space  $\langle \mathcal{X}, (\cdot, \cdot) \rangle$ , for any  $\mathbf{w} \in \mathcal{X}$ , we obtain a linear function  $f_{\mathbf{w}}$  from  $\mathcal{X}$  to  $\mathbf{R}$  by defining

$$f_{\mathbf{w}}(\mathbf{x}) = (\mathbf{w}, \mathbf{x}). \quad (1)$$

Throughout the paper, we define the (square) loss of prediction  $\hat{y}_t$  on the pair  $(x_t, y_t)$  by the squared prediction error  $(\hat{y}_t - y_t)^2$  and, accordingly, define the total loss of a sequence of predictions by the sum  $\sum_t (\hat{y}_t - y_t)^2$  of the squared prediction errors.

Typically, we express the bounds on the loss of our algorithms as a function of

$$\inf_{\mathbf{w}} \sum_t ((\mathbf{w}, \mathbf{x}_t) - y_t)^2, \quad (2)$$

<sup>1</sup>The general results will hold for finite and infinite dimensional vector spaces.

where the infimum is taken over all  $\mathbf{w}$  whose norm  $\sqrt{\langle \mathbf{w}, \mathbf{w} \rangle}$  is bounded by a parameter. Roughly speaking, this quantity measures the total misfit or noise of the environment with respect to the best “model” in the inner product space. In other words, bounds in terms of (2) are strong to the extent that there is a (not too large)  $\mathbf{w}$  for which  $f\mathbf{w}$  “approximately” maps  $\mathbf{x}_t$ ’s to corresponding  $y_t$ ’s. Thus we can also interpret (2) as an approximation error with respect to some unknown law generating the pairs  $(x_t, y_t)$ . In practice, estimates of this approximation error could be obtained by looking a posteriori at related applications, and calculating the quantity analogous to (2) for them.

In many cases we can even bound the additional loss of the algorithm over the above infimum similarly to the additional loss bounds of [3] obtained in a simpler setting. Our bounds are worst-case in the sense that they hold for all sequences of pairs  $(\mathbf{x}_t, y_t)$ . (In some cases we assume the norm of the  $\mathbf{x}_t$ ’s is bounded by a second parameter.)

Faber and Mycielski [6] noted that a natural class of smooth functions of a single real variable can be defined using inner products as above. The same class of smooth functions, as well as linear functions in  $\mathbf{R}^n$ , has been heavily studied in statistics [9] (however, with probabilistic assumptions). Thus, general results for learning classes of functions defined by arbitrary inner product spaces can be applied in a variety of circumstances. Faber and Mycielski proved bounds on  $\sum_t (\hat{y}_t - y_t)^2$  under the assumption that there was a  $\mathbf{w} \in \mathcal{X}$  for which for all  $t$ ,  $y_t = \langle \mathbf{w}, \mathbf{x}_t \rangle$ , and described some applications of this result for learning classes of smooth functions. Mycielski [24] had already treated the special case of linear functions in  $\mathbf{R}^n$ . The algorithm they analyzed for this “noise-free” case was a generalization of the on-line gradient descent algorithm<sup>2</sup> to arbitrary inner product spaces. We call this algorithm **GD** (defined below). In this paper we analyze the behavior of **GD** in the case in which there isn’t necessarily a  $\mathbf{w}$  for which for all  $t$ ,  $y_t = \langle \mathbf{w}, \mathbf{x}_t \rangle$ . Faber and Mycielski [6] also studied this case, but their algorithms made use of side information which, in this paper, we assume is not available. Hui and Zak [12] also studied the robustness of **GD** in the presence of noise in a similar setting, however they modelled observation noise, assuming that there was a  $\mathbf{w}$  such that for all  $t$   $y_t = \langle \mathbf{w}, \mathbf{x}_t \rangle$ , but that the learner’s observation of  $y_t$  was corrupted with noise. A more substantive difference is that they assumed  $\mathbf{x}_1 = \mathbf{x}_2 = \mathbf{x}_3 \dots$

The algorithm **GD** for the special case of linear functions in  $\mathbf{R}^n$  is a central building block in area of signal processing (see e.g. [19], [25], [26], [29], [10]) where it is usually called *Least-Mean-Square* (LMS) algorithm. Therefore, there is an extensive literature studying the convergence properties of this algorithm (see, e.g. [29], [10]). All this research, however, is based on probabilistic assumptions on the generation of the  $\mathbf{x}_t$ ’s. This paper shows that the algorithm **GD** can be analyzed even without probabilistic assumptions.

<sup>2</sup>Even though in the neural network community this algorithm is usually credited to Widrow and Hoff [28], a similar algorithm for the iterative solution of a system of linear equations was previously developed by Kaczmarz [13].

Gradient descent is an algorithm design technique which has achieved considerable practical success in more complicated hypothesis spaces, in particular multilayer neural networks. Despite this success, there appears not to be a principled method for tuning the learning rate. In this paper, we tune the learning rate in presence of noise with the goal of minimizing the worst-case total squared loss over the best that can be obtained using elements from a given class of linear functions.

The **GD** algorithm maintains an element  $\hat{\mathbf{w}}$  of  $\mathcal{X}$  as its hypothesis which is updated between trials. For each  $t$ , let  $\hat{\mathbf{w}}_t$  be the hypothesis before trial  $t$  (the initial hypothesis  $\hat{\mathbf{w}}_1$  is the zero vector). **GD** predicts with  $\hat{y}_t = \langle \hat{\mathbf{w}}_t, \mathbf{x}_t \rangle$  and updates the hypothesis following the rule

$$\hat{\mathbf{w}}_{t+1} = \hat{\mathbf{w}}_t - \eta(\hat{y}_t - y_t)\mathbf{x}_t. \quad (3)$$

where  $\eta > 0$  is the learning rate parameter.

If the real vector space  $\mathcal{X}$  has finite dimension, then each element  $\mathbf{v}$  of  $\mathcal{X}$  can be uniquely represented by the real vector  $\mathbf{c}(\mathbf{v})$  of its Fourier coefficients, once a basis is chosen. If the basis is orthonormal, by simple linear algebra facts we have  $\hat{y}_t = \langle \hat{\mathbf{w}}_t, \mathbf{x}_t \rangle = \mathbf{c}(\hat{\mathbf{w}}_t) \cdot \mathbf{c}(\mathbf{x}_t)$ . Furthermore, the vector  $2(\hat{y}_t - y_t)\mathbf{c}(\mathbf{x}_t)$  is the gradient, with respect to the vector  $\mathbf{c}(\hat{\mathbf{w}}_t)$ , of the square loss  $(\hat{y}_t - y_t)^2$  for the pair  $(\mathbf{x}_t, y_t)$ . Hence, in this case, rule (3) is indeed an “on-line” version of gradient descent performed over the quadratic loss.

When  $\mathcal{X}$  is an arbitrary real vector space, and therefore its elements may not be uniquely represented by finite tuples of reals, the **GD** algorithm is a natural generalization of on-line gradient descent<sup>3</sup> and may be viewed as follows [23].<sup>4</sup> After each trial  $t$ , there is a set  $S_t$  of elements  $\mathbf{w}$  of  $\mathcal{X}$  for which  $\langle \mathbf{w}, \mathbf{x}_t \rangle = y_t$ . Intuitively, our hypothesis would like to be more like the elements of  $S_t$ , since we are banking on there being a nearly functional relationship  $f\mathbf{w}$  between the  $\mathbf{x}_t$ ’s and the  $y_t$ ’s. It does not want to change *too* much, however, because the example  $(\mathbf{x}_t, y_t)$  may be misleading. The **GD** algorithm “takes a step” in the direction of the element of  $S_t$  which is closest to  $\hat{\mathbf{w}}_t$  (using the natural notion of the distance between elements of an inner product space).

## II. OVERVIEW OF RESULTS

We now give an overview of the bounds obtained in this paper. We will use  $\langle s_t \rangle_t$  to denote sequences  $s_1, s_2, \dots, s_t, \dots$ , and  $S^*$  to denote the set of all finite sequences (empty sequence included) over a set  $S$ .

For any  $\mathbf{v} \in \mathcal{X}$ ,  $\|\mathbf{v}\| = \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle}$  measures the “size” of  $\mathbf{v}$ . We show in Theorem IV.3 that for all sequences  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_t \in (\mathcal{X} \times \mathbf{R})^*$  and for all positive reals  $X, W$ ,

<sup>3</sup>To be precise, if  $\mathcal{X}$  has countably infinite dimension, then **GD** can still be viewed as a mapping performing on-line gradient descent. Such a mapping is clearly noncomputable in general since each step might involve the update of an infinite number of coefficients. However, note that the  $t$ -th hypothesis  $\hat{\mathbf{w}}_t$  is a linear combination of the first  $t - 1$  examples  $\{\mathbf{x}_1, \dots, \mathbf{x}_{t-1}\}$  and can thus be represented by  $t - 1$  real coefficients.

<sup>4</sup>Actually, this interpretation was shown only in the slightly more restricted case that  $\langle \mathcal{X}, \langle \cdot, \cdot \rangle \rangle$  is a Hilbert space.

and  $E$ , if  $\max_t \|\mathbf{x}_t\| \leq X$  and  $L_W(\mathbf{s}) \leq E$ , where

$$L_W(\mathbf{s}) = \inf_{\|\mathbf{w}\| \leq W} \sum_t ((\mathbf{w}, \mathbf{x}_t) - y_t)^2,$$

then the **GD** algorithm (with learning rate tuned to  $X, W$ , and  $E$ ) achieves the following

$$\sum_t (\hat{y}_t - y_t)^2 \leq L_W(\mathbf{s}) + 2(WX)\sqrt{E} + (WX)^2. \quad (4)$$

(Notice that  $L_W(\mathbf{s}) \geq L_{W'}(\mathbf{s})$  for all  $W' \geq W$ .) The above bound is tight in a very strong sense: We show in Theorem VII.1 a lower bound of  $L_W(\mathbf{s}) + 2(WX)\sqrt{E} + (WX)^2$  that holds for all  $X, W$ , and  $E$ , also when these parameters are given to the algorithm ahead of time.

We then remove the assumption that a bound  $E$  on  $L_W(\mathbf{s})$  is known for some  $W$ . However, we require that  $y_t$ 's are in a certain range  $[-Y, Y]$  for some  $Y > 0$ . In Theorem IV.4 we show that for all positive reals  $X$  and  $Y$  and for all sequences  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_t \in (\mathcal{X} \times [-Y, Y])^*$  such that  $\max_t \|\mathbf{x}_t\| \leq X$ , the total loss incurred on  $\mathbf{s}$  by a variant of the **GD** algorithm (with learning rate tuned to the remaining parameters  $X$  and  $Y$ ) is at most

$$L_{Y/X}(\mathbf{s}) + 9.2 \left( Y \sqrt{L_{Y/X}(\mathbf{s})} + Y^2 \right). \quad (5)$$

Notice that the above result also holds when  $L_{Y/X}(\mathbf{s})$  is replaced by  $L_W(\mathbf{s})$  for any  $W \leq Y/X$ . Observe that  $\sum_t (\hat{y}_t - y_t)^2 - L_{Y/X}(\mathbf{s})$  can be interpreted as the excess of the algorithm's total loss over the best that can be obtained using vectors  $\mathbf{w}$  whose norms are at most  $Y/X$ . The above bound is tight within constant factors: We show in Theorem VII.2 that for all prediction algorithms  $A$  and all  $X, Y, E > 0$ , there is a sequence  $\mathbf{s}$  on  $\mathcal{X} \times [-Y, Y]$  such that  $\max_t \|\mathbf{x}_t\| = X$ ,  $L_{Y/X}(\mathbf{s}) = E$ , and the total squared loss of  $A$  on  $\mathbf{s}$  is at least  $E + 2Y\sqrt{E} + Y^2$ . However, the dimension of the inner product space must increase as a function of  $E$ . As before, the lower bound holds also if all three parameters are given to the algorithm ahead of time.

We continue by giving the algorithm less information about the sequence. For the case when only a bound  $X$  on the norm of any  $\mathbf{x}_t$  is known, we show in Theorem IV.1 that the **GD** algorithm, tuned to  $X$ , achieves the following upper bound on the total loss (the sum of its squared prediction errors):

$$2.25 \inf_{\mathbf{w} \in \mathcal{X}} \left[ \left( \max_t \|\mathbf{x}_t\|^2 \right) \|\mathbf{w}\|^2 + \sum_t ((\mathbf{w}, \mathbf{x}_t) - y_t)^2 \right]$$

on any sequence  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_t \in (\mathcal{X} \times \mathbf{R})^*$  such that  $\max_t \|\mathbf{x}_t\| \leq X$ . Note that this result shows how the **GD** algorithm is able to trade-off between the “size” of a  $\mathbf{w}$ , represented by its norm, and the extent to which  $\mathbf{w}$  “fits” the data sequence, represented by the total loss incurred by  $f\mathbf{w}$ .

Finally, with no assumptions on the environment of the learner, a further variant of the **GD** algorithm has the

following bound on the total loss (Theorem IV.6)

$$9 \inf_{\mathbf{w} \in \mathcal{X}} \left[ \left( \max_t \|\mathbf{x}_t\|^2 \right) \|\mathbf{w}\|^2 + \sum_t ((\mathbf{w}, \mathbf{x}_t) - y_t)^2 \right]$$

that holds on any sequence  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_t \in (\mathcal{X} \times \mathbf{R})^*$ .

Our general results are summarized in Table I.

We may apply our general bounds to a class of smooth functions of a single real variable, in the manner used by Faber and Mycielski [6] in the case that there is a perfect smooth function. The smoothness of a function is measured by the 2-norm of its derivative. Of course, the derivative measures the steepness of a function at a given point, and therefore the 2-norm (or any norm, for that matter) of the derivative measures the tendency of the function to be steep. When normalized appropriately, the 2-norm of a function  $f$ 's derivative can be seen to be between the average steepness of  $f$  and the  $f$ 's maximum steepness. In Theorem V.1 we show that if there is an absolutely continuous function  $f : \mathbf{R}_+ \rightarrow \mathbf{R}$  with  $f(0) = 0$  which tends not to be very steep and approximately maps the  $x_t$ 's to the  $y_t$ 's, and if the  $x_t$ 's are not very big, then an application of the **GD** algorithm to this case obtains good bounds on the total loss. More formally, we show that, for example, if the  $x_t$ 's are taken from  $[0, X]$ , and if  $f : [0, \infty) \rightarrow \mathbf{R}$  satisfies  $\|f'\|_2 = \sqrt{\int_0^X f'(u)^2 du} \leq W$ , and  $\sum_t (f(x_t) - y_t)^2 \leq E$ , then the special case of the general **GD** algorithm applied to this problem has a sum of squared errors bounded by

$$\inf_{\|f'\|_2 \leq W} \left[ \sum_t (f(x_t) - y_t)^2 \right] + 2W\sqrt{XE} + W^2X. \quad (6)$$

A bound of  $W^2X$  was proved by [6] in the case when  $E = 0$ . It is surprising that the time required for the algorithm we describe for this problem to make its  $t$ th prediction  $\hat{y}_t$  is  $O(t)$  in the uniform cost model provided that all past examples and predictions are saved. This is because, although the vector space in which we live in this application consists of functions, and therefore the **GD** algorithm requires us to add functions, we can see that the functions that arise are piecewise linear, with the pieces being a simple functions of the past examples and predictions. In the case  $E = 0$ , however, there is an algorithm with an optimal bound on  $\sum_t (\hat{y}_t - y_t)^2$  which computes its  $t$ th prediction in  $O(\log t)$  time [14], raising the hope that there might be a similarly efficient robust algorithm. In Theorem V.2 we extend our result to apply to classes of smooth functions of  $n > 1$  real variables studied by Faber and Mycielski [6] in the absence of noise. We further show that upper bound (6), even viewed as bound on the excess of the algorithm's total loss over the loss of the best function of “size” at most  $W$ , is optimal, constants included.

Littlestone, Long and, Warmuth [18] proved bounds for another algorithm for learning linear functions in  $\mathbf{R}^n$ , in which the  $\mathbf{x}_t$ 's were measured using the infinity norm, and the  $\mathbf{w}$ 's were measured using 1-norm. The bounds for the two algorithms are incomparable because different norms

Known information	Bounds for $L(\mathbf{s}) = \sum_t (\hat{y}_t - y_t)^2$
$\max_t \ \mathbf{x}_t\  \leq X$ $L_W(\mathbf{s}) \leq E$	$L(\mathbf{s}) \leq L_W(\mathbf{s}) + 2(WX)\sqrt{E} + (WX)^2$ MATCHING LOWER BOUND
$\max_t \ \mathbf{x}_t\  \leq X$ $\max_t  y_t  \leq Y$	$L(\mathbf{s}) \leq L_{Y/X}(\mathbf{s}) + 9.2(Y\sqrt{L_{Y/X}(\mathbf{s})} + Y^2)$ $L(\mathbf{s}') \geq L_{Y/X}(\mathbf{s}') + 2Y\sqrt{L_{Y/X}(\mathbf{s}')} + Y^2$
$\max_t \ \mathbf{x}_t\  \leq X$	$L(\mathbf{s}) \leq 2.25 \inf_{\mathbf{w} \in \mathcal{X}} [(\max_t \ \mathbf{x}_t\ ^2) \ \mathbf{w}\ ^2 + \sum_t ((\mathbf{w}, \mathbf{x}_t) - y_t)^2]$
None	$L(\mathbf{s}) \leq 9 \inf_{\mathbf{w} \in \mathcal{X}} [(\max_t \ \mathbf{x}_t\ ^2) \ \mathbf{w}\ ^2 + \sum_t ((\mathbf{w}, \mathbf{x}_t) - y_t)^2]$

TABLE I

TABLE SUMMARIZING LOSS BOUNDS FOR **GD** AND ITS VARIANTS WITH RESPECT TO DIFFERENT A PRIORI INFORMATION.

are used to measure the sizes of the  $\mathbf{x}$ 's and the  $\mathbf{w}$ 's. However, the algorithm of [18] does not appear to generalize to arbitrary inner product spaces as did the **GD** algorithm, and therefore those techniques do not appear to be as widely applicable.

One of the main problems with gradient descent is that it motivates a learning rule but does not give any method for choosing the step size. Our results provide a method for setting the learning rate essentially optimally when learning linear functions. An exciting research direction is to investigate to what extent the methods of this paper can be applied to analyze other simple gradient descent learning algorithms.

Our methods can also be applied to the batch setting where the whole sequence of examples is given to the learner at once and the goal of learning is to find the function that minimizes the sum of the squared residual errors. In the case of linear functions this can be solved directly using the linear least squares method which might be considered to be too computationally expensive. Iterative methods provide an alternative. We prove a total loss bound for a gradient descent algorithm by applying the techniques used in this paper. We then contrast this bound to the standard bound for steepest descent on the total squared residual error.

The paper is organized as follows: In Section III we recall the notion of inner product space and define the algorithm **GD**. The upper bounds for **GD** and its variants are all proven in Section IV; in this section we also prove bounds for the normalized total loss. These results are applied in Section V to derive upper bounds for prediction in classes of smooth functions. The comparison with the standard steepest descent methods is given in Section VI.

Corresponding lower bounds for the upper bounds of Sections IV and V are then proven in Section VII. The paper is concluded in Section VIII with some discussion and open problems.

### III. PRELIMINARIES

Let  $\mathbf{N}$  denote the positive integers,  $\mathbf{R}$  denote the reals. Each prediction of an on-line algorithm is determined by the previous examples and the current instance. In this paper the domain of the instances is always a fixed real vector space  $\mathcal{X}$ . An on-line prediction algorithm  $A$  is a mapping from  $(\mathcal{X} \times \mathbf{R})^* \times \mathcal{X}$  to  $\mathbf{R}$ . For a finite sequence  $\mathbf{s} = \langle (x_t, y_t) \rangle_{1 \leq t \leq m}$  of examples we let  $\hat{y}_t$  denote the *prediction* of  $A$  on the  $t$ -th *trial*, i.e.,

$$\hat{y}_t = A(\langle (x_1, y_1), \dots, (x_{t-1}, y_{t-1}) \rangle, x_t).$$

and we call  $\hat{y}_1, \dots, \hat{y}_m$  the *sequence of  $A$ 's on-line predictions for  $\mathbf{s}$* .

An *inner product space* (sometimes called a pre-Hilbert space since the imposition of one more assumption yields the definition of a Hilbert space) consists of a real vector space  $\mathcal{X}$  and a function  $(\cdot, \cdot)$  (called an *inner product*) from  $\mathcal{X} \times \mathcal{X}$  to  $\mathbf{R}$  that satisfies the following for all  $\mathbf{u}, \mathbf{v}, \mathbf{x} \in \mathcal{X}$  and  $\kappa \in \mathbf{R}$ :

1.  $(\mathbf{u}, \mathbf{v}) = (\mathbf{v}, \mathbf{u})$ ;
2.  $(\kappa \mathbf{u}, \mathbf{v}) = \kappa(\mathbf{u}, \mathbf{v})$ ;
3.  $(\mathbf{u} + \mathbf{v}, \mathbf{x}) = (\mathbf{u}, \mathbf{x}) + (\mathbf{v}, \mathbf{x})$ ;
4.  $(\mathbf{x}, \mathbf{x}) > 0$  whenever  $\mathbf{x} \neq \mathbf{0}$ .

The last requirement can be dropped essentially without affecting the definition (see e.g. [30, page 25]). For  $\mathbf{x} \in \mathcal{X}$ , the *norm* of  $\mathbf{x}$ , denoted by  $\|\mathbf{x}\|$ , is defined by

$$\|\mathbf{x}\| = \sqrt{(\mathbf{x}, \mathbf{x})}.$$

**Algorithm GD.****Input:**  $\eta \geq 0$ .

- Choose  $\mathcal{X}$ 's zero vector as initial hypothesis  $\hat{\mathbf{w}}_1$ .
- On each trial  $t$ :
  1. Get  $\mathbf{x}_t \in \mathcal{X}$  from the environment.
  2. Predict with  $\hat{y}_t = (\hat{\mathbf{w}}_t, \mathbf{x}_t)$ .
  3. Get  $y_t \in \mathcal{Y}$  from the environment.
  4. Update the current hypothesis  $\hat{\mathbf{w}}_t$  according to the rule

$$\hat{\mathbf{w}}_{t+1} = \hat{\mathbf{w}}_t + \eta(y_t - \hat{y}_t)\mathbf{x}_t.$$

Fig. 1. Pseudo-code for algorithm **GD**. (See Theorems IV.1, IV.2, IV.3, and Corollary IV.1.)

(These definitions are taken from [30].)

An example of an inner product is the dot product in  $\mathbf{R}^n$ . For  $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$  for some positive integer  $n$ , the dot product of  $\mathbf{x}$  and  $\mathbf{y}$  is defined to be

$$\mathbf{x} \cdot \mathbf{y} = \sum_{i=1}^n x_i y_i.$$

The 2-norm (or Euclidian norm) of  $\mathbf{x} \in \mathbf{R}^n$  is then defined to be

$$\|\mathbf{x}\|_2 = \sqrt{\mathbf{x} \cdot \mathbf{x}} = \sqrt{\sum_{i=1}^n x_i^2}.$$

If  $f$  is a function from  $\mathbf{R}$  to  $\mathbf{R}$ , we say that  $f$  is *absolutely continuous*<sup>5</sup> iff there exists a (Lebesgue measurable) function  $g : \mathbf{R} \rightarrow \mathbf{R}$  such that for all  $a, b \in \mathbf{R}$ ,  $a \leq b$ ,

$$f(b) - f(a) = \int_a^b g(x) dx.$$

#### IV. UPPER BOUNDS FOR THE GENERALIZED GRADIENT DESCENT ALGORITHM

In this section, we prove bounds on the worst case total loss made by the **GD** algorithm (described in Figure 1). (Technically, Figure 1 describes a different learning algorithm for each initial setting of the “learning rate”  $\eta$ . For a particular  $\eta$ , we will refer to the associated learning algorithm as **GD** $_\eta$ , and we will use a similar convention throughout the paper).

For the remainder of this section, fix an inner product space  $\langle \mathcal{X}, (\cdot, \cdot) \rangle$ . In what follows, we will analyze the **GD** algorithm and its variants starting from the case where only a bound on the norm of  $\mathbf{x}_t$ , for all  $t$ , is available to the learner ahead of time. We will then show how additional information can be exploited for tuning the learning rate  $\eta$  and obtaining better worst-case bounds. Finally, we will prove a bound for the case where no assumptions are made on the environment of the learner.

##### A. Bounding the size of the instances

In this section we prove that, when given a bound on  $\max_t \|\mathbf{x}_t\|$ , the algorithm **GD** can obtain good bounds on

<sup>5</sup>This is shown to be equivalent to a more technical definition in most Calculus texts.

the total loss. We will remove the assumption of this knowledge later through application of standard doubling techniques.

As a first step, we will show the following which might be interpreted as determining the “progress” per trial, that is the amount that **GD** $_\eta$  learns from an error. The derivation is based on previous derivations used in the proof of convergence of the on-line gradient descent algorithm (see, e.g. [5]).

*Lemma IV.1:* Choose  $\mathbf{x}, \hat{\mathbf{w}}_1, \mathbf{w} \in \mathcal{X}, y \in \mathbf{R}, \eta > 0$ . Let  $\hat{y} = (\hat{\mathbf{w}}_1, \mathbf{x})$  and  $\hat{\mathbf{w}}_2 = \hat{\mathbf{w}}_1 + \eta(y - \hat{y})\mathbf{x}$ . Then

$$\begin{aligned} & \|\hat{\mathbf{w}}_1 - \mathbf{w}\|^2 - \|\hat{\mathbf{w}}_2 - \mathbf{w}\|^2 \\ &= (2\eta - \eta^2 \|\mathbf{x}\|^2)(\hat{y} - y)^2 - 2\eta(y - \hat{y})(y - (\mathbf{w}, \mathbf{x})). \end{aligned} \quad (7)$$

**Proof:** Let  $\alpha = \eta(y - \hat{y})$ . Then  $\hat{\mathbf{w}}_2 = \hat{\mathbf{w}}_1 + \alpha\mathbf{x}$ . Thus

$$\begin{aligned} & \|\hat{\mathbf{w}}_2 - \mathbf{w}\|^2 \\ &= ((\hat{\mathbf{w}}_2 - \mathbf{w}), (\hat{\mathbf{w}}_2 - \mathbf{w})) \\ &= ((\hat{\mathbf{w}}_1 + \alpha\mathbf{x} - \mathbf{w}), (\hat{\mathbf{w}}_1 + \alpha\mathbf{x} - \mathbf{w})) \\ &= \|\hat{\mathbf{w}}_1 - \mathbf{w}\|^2 + (2\alpha\mathbf{x}, (\hat{\mathbf{w}}_1 - \mathbf{w})) + \alpha^2 \|\mathbf{x}\|^2. \end{aligned}$$

This implies

$$\begin{aligned} & \|\hat{\mathbf{w}}_2 - \mathbf{w}\|^2 - \|\hat{\mathbf{w}}_1 - \mathbf{w}\|^2 \\ &= 2\alpha(\mathbf{x}, (\hat{\mathbf{w}}_1 - \mathbf{w})) + \alpha^2 \|\mathbf{x}\|^2 \\ &= 2\alpha(\hat{y} - (\mathbf{w}, \mathbf{x})) + \alpha^2 \|\mathbf{x}\|^2 \\ &= 2\alpha(\hat{y} - y) + 2\alpha(y - (\mathbf{w}, \mathbf{x})) + \alpha^2 \|\mathbf{x}\|^2. \end{aligned}$$

Expanding our definition of  $\alpha$ ,

$$\begin{aligned} & \|\hat{\mathbf{w}}_2 - \mathbf{w}\|^2 - \|\hat{\mathbf{w}}_1 - \mathbf{w}\|^2 \\ &= -2\eta(\hat{y} - y)^2 + 2\eta(y - \hat{y})(y - (\mathbf{w}, \mathbf{x})) \\ & \quad + \eta^2 \|\mathbf{x}\|^2 (y - \hat{y})^2 \\ &= -(2\eta - \eta^2 \|\mathbf{x}\|^2)(\hat{y} - y)^2 + 2\eta(y - \hat{y})(y - (\mathbf{w}, \mathbf{x})), \end{aligned}$$

establishing (7).  $\square$

We need the following simple lemma:

*Lemma IV.2:* For all  $q, r, c \in \mathbf{R}$  such that  $0 < c < 1$ ,

$$q^2 - qr \geq cq^2 - \frac{r^2}{4(1-c)}. \quad (8)$$

**Proof.** Inequality (8) is equivalent to

$$\frac{(2(1-c)q - r)^2}{4(1-c)} \geq 0$$

which is clearly true as  $0 < c < 1$ .  $\square$

As a second step, we show a lower bound on the progress per trial. This lower bound will be used to prove the main theorem of this section.

*Lemma IV.3:* Choose  $\mathbf{x}, \hat{\mathbf{w}}_1, \mathbf{w} \in \mathcal{X}, y \in \mathbf{R}$ . Choose  $X, \beta, c \in \mathbf{R}$  such that  $X \geq \|\mathbf{x}\|$ ,  $0 < \beta < 2$  and  $0 < c < 1$ . Let

$$\hat{y} = (\hat{\mathbf{w}}_1, \mathbf{x}) \quad \text{and} \quad \hat{\mathbf{w}}_2 = \hat{\mathbf{w}}_1 + \frac{\beta}{X^2}(y - \hat{y})\mathbf{x}.$$

Then

$$\begin{aligned} & \|\hat{\mathbf{w}}_1 - \mathbf{w}\|^2 - \|\hat{\mathbf{w}}_2 - \mathbf{w}\|^2 \\ & \geq \frac{2\beta - \beta^2}{X^2} \left[ c(\hat{y} - y)^2 - \frac{\beta^2}{(2\beta - \beta^2)^2(1-c)}(y - (\mathbf{w}, \mathbf{x}))^2 \right]. \end{aligned}$$

**Proof.** Applying Lemma IV.1 with  $\eta = \frac{\beta}{X^2}$ , we get

$$\begin{aligned} & \|\hat{\mathbf{w}}_1 - \mathbf{w}\|^2 - \|\hat{\mathbf{w}}_2 - \mathbf{w}\|^2 \\ &= \left( \frac{2\beta}{X^2} - \frac{\beta^2 \|\mathbf{x}\|^2}{X^4} \right) (\hat{y} - y)^2 \\ &\quad - \frac{2\beta}{X^2} (y - \hat{y})(y - (\mathbf{w}, \mathbf{x})) \end{aligned} \quad (9)$$

$$\begin{aligned} &\geq \left( \frac{2\beta}{X^2} - \frac{\beta^2}{X^2} \right) (\hat{y} - y)^2 \\ &\quad - \frac{2\beta}{X^2} (y - \hat{y})(y - (\mathbf{w}, \mathbf{x})) \\ &\geq \frac{2\beta - \beta^2}{X^2} [(\hat{y} - y)^2 \\ &\quad - \frac{2\beta}{2\beta - \beta^2} |y - \hat{y}| |y - (\mathbf{w}, \mathbf{x})|] \\ &\geq \frac{2\beta - \beta^2}{X^2} [c(\hat{y} - y)^2 \\ &\quad - \frac{\beta^2}{(2\beta - \beta^2)^2(1-c)} (y - (\mathbf{w}, \mathbf{x}))^2] \end{aligned} \quad (10)$$

where Inequality (10) holds because  $X \geq \|\mathbf{x}\|$  and Inequality (11) is an application of Lemma IV.2.  $\square$

The next theorem shows that the performance of the **GD** algorithm degrades gracefully as the relationship to be modelled moves away from being  $(\mathbf{w}, \cdot)$  from some  $\mathbf{w} \in \mathcal{X}$ . Throughout the paper, for all sequences  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_t \in (\mathcal{X} \times \mathbf{R})^*$  and all  $\mathbf{w} \in \mathcal{X}$ , let

$$L_{\mathbf{w}}(\mathbf{s}) = \sum_t ((\mathbf{w}, \mathbf{x}_t) - y_t)^2,$$

and for all  $W > 0$  let

$$L_W(\mathbf{s}) = \inf_{\|\mathbf{w}\| \leq W} L_{\mathbf{w}}(\mathbf{s}).$$

*Theorem IV.1:* Choose  $0 < \beta < 2$ ,  $0 < c < 1$ ,  $m \in \mathbf{N}$ , and  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_{t \leq m} \in (\mathcal{X} \times \mathbf{R})^m$ . Let  $X \geq \max_t \|\mathbf{x}_t\|$ , and let  $\hat{y}_1, \dots, \hat{y}_m$  be the sequence of  $\mathbf{GD}_{\beta/X^2}$ 's on-line predictions for  $\mathbf{s}$ . Then,

$$\sum_{t=1}^m (\hat{y}_t - y_t)^2 \leq \inf_{\mathbf{w} \in \mathcal{X}} \left[ \frac{X^2 \|\mathbf{w}\|^2}{(2\beta - \beta^2)c} + \frac{L_{\mathbf{w}}(\mathbf{s})}{(2 - \beta)^2 c(1 - c)} \right]. \quad (12)$$

In particular, if  $\beta = 2/3$  and  $c = 1/2$ ,

$$\sum_{t=1}^m (\hat{y}_t - y_t)^2 \leq 2.25 \inf_{\mathbf{w} \in \mathcal{X}} [X^2 \|\mathbf{w}\|^2 + L_{\mathbf{w}}(\mathbf{s})]. \quad (13)$$

Notice that, by setting  $c = 1/2$  and by letting  $\beta \rightarrow 0$ , the constant on the  $L_{\mathbf{w}}(\mathbf{s})$  term can be brought arbitrarily close to 1 at the expense of increasing the constant on the other term.

**Proof:** Choose  $\mathbf{w} \in \mathcal{X}$ . If  $\hat{\mathbf{w}}_1, \hat{\mathbf{w}}_2, \dots, \hat{\mathbf{w}}_{m+1}$  is the sequence of  $\mathbf{GD}_{\beta/X^2}$ 's hypotheses, we get

$$\begin{aligned} & \sum_{t=1}^m \frac{2\beta - \beta^2}{X^2} \left[ c(\hat{y}_t - y_t)^2 - \frac{\beta^2}{(2\beta - \beta^2)^2(1-c)} (y_t - (\mathbf{w}, \mathbf{x}_t))^2 \right] \\ & \leq \sum_{t=1}^m (\|\hat{\mathbf{w}}_t - \mathbf{w}\|^2 - \|\hat{\mathbf{w}}_{t+1} - \mathbf{w}\|^2) \quad \text{by Lemma IV.3} \\ & = \|\hat{\mathbf{w}}_1 - \mathbf{w}\|^2 - \|\hat{\mathbf{w}}_{m+1} - \mathbf{w}\|^2 \\ & \leq \|\mathbf{w}\|^2 \quad \text{since } \hat{\mathbf{w}}_1 = \vec{0} \text{ and } \|\cdot\| \text{ is nonnegative.} \end{aligned}$$

Thus

$$\begin{aligned} & \sum_{t=1}^m \left[ c(\hat{y}_t - y_t)^2 - \frac{\beta^2}{(2\beta - \beta^2)^2(1-c)} (y_t - (\mathbf{w}, \mathbf{x}_t))^2 \right] \\ & \leq \frac{X^2 \|\mathbf{w}\|^2}{2\beta - \beta^2}. \end{aligned}$$

Solving for  $\sum_t (\hat{y}_t - y_t)^2$  yields

$$\sum_{t=1}^m (\hat{y}_t - y_t)^2 \leq \frac{X^2 \|\mathbf{w}\|^2}{(2\beta - \beta^2)c} + \frac{\beta^2}{(2\beta - \beta^2)^2 c(1-c)} L_{\mathbf{w}}(\mathbf{s})$$

establishing (12). Formula (13) then follows immediately.  $\square$

Observe that the assumption  $\hat{\mathbf{w}}_1 = \vec{0}$  is chosen merely for convenience. If  $\hat{\mathbf{w}}_1 \neq \vec{0}$ , then the factor  $\|\mathbf{w}\|^2$  in (12) is replaced by  $\|\mathbf{w} - \hat{\mathbf{w}}_1\|^2$ . Thus, in this more general form, the bound of Theorem IV.1 depends on the squared distance between the starting vector  $\hat{\mathbf{w}}_1$  and the ‘‘target’’  $\mathbf{w}$ .

### A.1 Normalized loss

If we run algorithm **GD** with learning rate  $\eta$  set in each trial  $t$  to  $\frac{\beta}{\|\mathbf{x}_t\|^2}$ , we can then prove a variant of Theorem IV.1 for a different notion of loss (previously studied by Faber and Mycielski [6]) which we call *normalized loss*. The normalized loss incurred by an algorithm predicting  $\hat{y}_t$  on a trial  $(\mathbf{x}_t, y_t)$  is defined by  $\frac{(\hat{y}_t - y_t)^2}{\|\mathbf{x}_t\|^2}$ . We begin by proving the following result via a straightforward variant of the proof of Lemma IV.3.

*Lemma IV.4:* Choose  $\mathbf{x}, \hat{\mathbf{w}}_1, \mathbf{w} \in \mathcal{X}, y \in \mathbf{R}, 0 < \beta < 2$ , and  $0 < c < 1$ . Let

$$\hat{y} = (\hat{\mathbf{w}}_1, \mathbf{x}) \quad \text{and} \quad \hat{\mathbf{w}}_2 = \hat{\mathbf{w}}_1 + \frac{\beta}{\|\mathbf{x}\|^2} (y - \hat{y}) \mathbf{x}.$$

Then

$$\begin{aligned} & \|\hat{\mathbf{w}}_1 - \mathbf{w}\|^2 - \|\hat{\mathbf{w}}_2 - \mathbf{w}\|^2 \\ & \geq \frac{2\beta - \beta^2}{\|\mathbf{x}\|^2} \left[ (\hat{y} - y)^2 c - \frac{\beta^2}{(2\beta - \beta^2)^2(1-c)} (y - (\mathbf{w}, \mathbf{x}))^2 \right]. \end{aligned}$$

We now extend Theorem IV.1 to the normalized loss. Let

$$L'_{\mathbf{w}}(\mathbf{s}) = \sum_{t=1}^m \frac{(f_{\mathbf{w}}(\mathbf{x}_t) - y_t)^2}{\|\mathbf{x}_t\|^2}.$$

*Theorem IV.2:* Choose  $0 < \beta < 2$ ,  $m \in \mathbf{N}$ , and  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_{t \leq m} \in (\mathcal{X} \times \mathbf{R})^m$ . Let  $\hat{y}_1, \dots, \hat{y}_m$  be the sequence of  $\mathbf{GD}_{\beta/\|\mathbf{x}_t\|^2}$ 's on-line predictions for  $\mathbf{s}$ . Then,

$$\sum_{t=1}^m \frac{(\hat{y}_t - y_t)^2}{\|\mathbf{x}_t\|^2} \leq \inf_{\mathbf{w} \in \mathcal{X}} \left[ \frac{\|\mathbf{w}\|^2}{\beta(2 - \beta)c} + \frac{L'_{\mathbf{w}}(\mathbf{s})}{(2 - \beta)^2 c(1 - c)} \right]$$

for all  $0 < c < 1$ . In particular, if  $\beta = 2/3$  and  $c = \frac{1}{2}$ ,

$$\sum_{t=1}^m \frac{(\hat{y}_t - y_t)^2}{\|\mathbf{x}_t\|^2} \leq 2.25 \inf_{\mathbf{w} \in \mathcal{X}} [\|\mathbf{w}\|^2 + L'_{\mathbf{w}}(\mathbf{s})].$$

The above theorem shows that the knowledge of a bound on  $\|\mathbf{x}_t\|$ , for all  $t$ , is not necessary when the normalized loss is used. This raises the question of whether the setting

$\eta = \frac{\beta}{\|\mathbf{x}_t\|^2}$  (for some fixed  $\beta$  not depending on  $\|\mathbf{x}_t\|$ ) can be successfully used when the goal is to minimize the total unnormalized loss and no bound on  $\|\mathbf{x}_t\|$  is available beforehand. On the other hand, suppose  $\mathcal{X} = \mathbf{R}$ , and the inner product is just the ordinary product on the reals. Suppose further that for  $\epsilon > 0$ ,  $x_1 = \epsilon$ , and  $y_1 = 1$ , whereas for all  $t > 1$ ,  $x_t = 1$  and  $y_t = 0$ . Then for smaller and smaller  $\epsilon$ , the total (unnormalized) quadratic loss of the **GD** with the above setting of  $\eta$  in this case is unbounded, whereas there is a  $w$  such that  $\sum_t (wx_t - y_t)^2 = 1$ , namely 0. (This example is due to Ethan Bernstein.)

### B. Tuning $\beta$

The next result shows that, if certain parameters are known in advance, optimal performance can be obtained by tuning  $\beta$ . We need a technical lemma first. Define the function  $G : \mathbf{R}_+^3 \rightarrow (0, 1]$  by

$$G(E, W, X) = \frac{WX}{\sqrt{E} + WX}.$$

*Lemma IV.5:* For all  $E, W, X > 0$

$$\frac{(WX)^2}{\beta(2-\beta)c} + \frac{E}{(2-\beta)^2c(1-c)} = E + (WX)^2 + 2WX\sqrt{E} \quad (14)$$

whenever  $\beta = G(E, W, X)$  and  $c = \frac{\sqrt{E}+WX}{2\sqrt{E}+WX}$ .

**Proof.** First notice that, when  $\beta$  and  $c$  are chosen as in the lemma's hypothesis,  $0 < \beta \leq 1$  and  $\frac{1}{2} \leq c < 1$  for all  $E, W, X \geq 0$ . Second, observe that (14) can be rewritten as

$$\frac{1}{\beta(2-\beta)c} + \frac{y^2}{(2-\beta)^2c(1-c)} = (y+1)^2 \quad (15)$$

where  $y = \frac{\sqrt{E}}{WX}$ . Now let

$$\beta = G(E, W, X) = \frac{1}{y+1} \quad \text{and} \quad c = \frac{y+1}{2y+1}.$$

Then

$$(2-\beta)c = 1 \quad \text{and} \quad (2-\beta)^2c(1-c) = 1 - \beta = \frac{y}{y+1}.$$

By making these substitutions in (15) we obtain  $y+1 + y(y+1) = (y+1)^2$ .  $\square$

*Theorem IV.3:* For each  $E, X, W \geq 0$ , the algorithm  $\mathbf{GD}_{G(E,W,X)/X^2}$  has the following properties. Choose  $m \in \mathbf{N}$ ,  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_{t \leq m} \in (\mathcal{X} \times \mathbf{R})^m$ , such that  $\max_t \|\mathbf{x}_t\| \leq X$ , and  $L_W(\mathbf{s}) \leq E$ . Let  $\hat{y}_1, \dots, \hat{y}_m$  be the sequence of  $\mathbf{GD}_{G(E,W,X)/X^2}$ 's on-line predictions for  $\mathbf{s}$ . Then,

$$\sum_{t=1}^m (\hat{y}_t - y_t)^2 \leq L_W(\mathbf{s}) + 2WX\sqrt{E} + (WX)^2.$$

**Proof.** Choose  $m \in \mathbf{N}$ ,  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_{t \leq m} \in (\mathcal{X} \times \mathbf{R})^m$  for which  $L_W(\mathbf{s}) \leq E$  and  $\max_t \|\mathbf{x}_t\|^2 \leq X$ . By Theorem IV.1,

for all  $\beta$  and  $c$  such that  $0 < \beta < 2$  and  $0 < c < 1$ , we have

$$\begin{aligned} & \sum_{t=1}^m (\hat{y}_t - y_t)^2 \\ & \leq \inf_{\mathbf{w} \in \mathcal{X}} \left[ \frac{X^2 \|\mathbf{w}\|^2}{\beta(2-\beta)c} + \frac{L_W(\mathbf{s})}{(2-\beta)^2c(1-c)} \right] \\ & \leq \frac{(WX)^2}{\beta(2-\beta)c} + \frac{L_W(\mathbf{s})}{(2-\beta)^2c(1-c)} \\ & = \frac{(WX)^2}{\beta(2-\beta)c} + \frac{L_W(\mathbf{s})}{(2-\beta)^2c(1-c)} - L_W(\mathbf{s}) + L_W(\mathbf{s}) \\ & \leq \frac{(WX)^2}{\beta(2-\beta)c} + \frac{E}{(2-\beta)^2c(1-c)} - E + L_W(\mathbf{s}) \end{aligned}$$

since  $L_W(\mathbf{s}) \leq E$  and  $0 < (2-\beta)^2c(1-c) \leq 1$  for the given ranges of  $c$  and  $\beta$ . Applying Lemma IV.5 for  $\beta = G(E, W, X)$  and  $c = \frac{\sqrt{E}+WX}{2\sqrt{E}+WX}$ , we then conclude

$$\begin{aligned} \sum_{t=1}^m (\hat{y}_t - y_t)^2 & \leq E + 2WX\sqrt{E} + (WX)^2 - E + L_W(\mathbf{s}) \\ & = L_W(\mathbf{s}) + 2WX\sqrt{E} + (WX)^2 \end{aligned}$$

as desired.  $\square$

A corollary to Theorem IV.3 can be obtained for the normalized loss.

*Corollary IV.1:* For any  $E, W \geq 0$  and for any  $m \in \mathbf{N}$ . Choose  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_{t \leq m} \in (\mathcal{X} \times \mathbf{R})^m$ , such that  $L'_W(\mathbf{s}) \leq E$ . Let  $\hat{y}_1, \dots, \hat{y}_m$  be the sequence of  $\mathbf{GD}_{G(E,W,1)/\|\mathbf{x}_t\|^2}$ 's on-line predictions for  $\mathbf{s}$ . Then,

$$\sum_{t=1}^m \frac{(\hat{y}_t - y_t)^2}{\|\mathbf{x}_t\|^2} \leq L'_W(\mathbf{s}) + 2W\sqrt{E} + W^2.$$

**Proof.** Lemma IV.5 can be applied to Theorem IV.2 with  $X = 1$  and  $c = \frac{\sqrt{E}+W}{2\sqrt{E}+W}$ . The derivation then closely resembles that of Theorem IV.3.  $\square$

The following corollary shows an application of Theorem IV.3 to the class of linear functions in  $\mathbf{R}^n$ . For each  $W, n$ , let  $\text{LIN}_{W,n}$  be the set of all functions  $f$  from  $\mathbf{R}^n$  to  $\mathbf{R}$  for which there exists  $\mathbf{w} \in \mathbf{R}^n$ ,  $\|\mathbf{w}\|_2 \leq W$ , such that for all  $\mathbf{x} \in \mathbf{R}^n$ ,  $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$ .

*Corollary IV.2:* For each  $E, X, W \geq 0$ , the **GD** algorithm has the following properties. Choose  $m, n \in \mathbf{N}$ ,  $\langle (\mathbf{x}_t, y_t) \rangle_{t \leq m} \in (\mathbf{R}^n \times \mathbf{R})^m$ , such that  $\max_t \|\mathbf{x}_t\|_2 \leq X$ , and there is an  $f \in \text{LIN}_{W,n}$  for which  $\sum_{t=1}^m (f(\mathbf{x}_t) - y_t)^2 \leq E$ . Let  $\hat{y}_1, \dots, \hat{y}_m$  be the sequence of  $\mathbf{GD}_{G(E,W,X)/X^2}$ 's on-line predictions for  $\mathbf{s}$ , when **GD** is applied to the inner product space  $\text{LIN}_{W,n}$ . Then,

$$\begin{aligned} & \sum_{t=1}^m (\hat{y}_t - y_t)^2 \\ & \leq \inf_{f \in \text{LIN}_{W,n}} \left[ \sum_{t=1}^m (f(\mathbf{x}_t) - y_t)^2 \right] + 2WX\sqrt{E} + W^2X^2. \end{aligned}$$

It has been shown recently [15] that even on very simple probabilistic artificial data, the above tunings and worst-case loss bounds are close to optimal.

In the next section, we show that techniques from [3] may also be applied to obtain a  $L_{Y/X}(\mathbf{s}) + O(Y\sqrt{E} + Y^2)$  bound on the total loss (unnormalized) when bounds  $X$  on  $\|\mathbf{x}_t\|$  and  $Y$  on  $|y_t|$  are known for all  $t$ . However, the delicate interplay between  $L_W$  and  $W$  (loosely speaking, increasing  $W$  decreases  $L_W$ ) has so far prevented us from obtaining such a result without knowledge of any of the three parameters  $W$ ,  $X$ , and  $E$ .

**Algorithm G1.****Input**  $X, Y \geq 0$ .

- For each  $i = 0, 1, \dots$ 
  - Let  $k_i = z^i (aY)^2$ .
  - Repeat
    1. Give  $\mathbf{x}_t$  to  $\mathbf{GD}_{G(k_i, Y/X, X)/X^2}$ .
    2. Get  $\mathbf{GD}_{G(k_i, Y/X, X)/X^2}$ 's prediction  $h_t$ .
    3. Predict with

$$\hat{y}_t = \begin{cases} -Y & \text{if } h_t < -Y \\ h_t & \text{if } |h_t| \leq Y \\ Y & \text{otherwise.} \end{cases}$$

4. Pass  $y_t$  to  $\mathbf{GD}_{G(k_i, Y/X, X)/X^2}$  until the total loss in this loop exceeds

$$k_i + 2Y\sqrt{k_i} + Y^2.$$

Fig. 2. Pseudo-code for the algorithm **G1**. (See Theorem IV.4.) Here  $\mathbf{GD}$  is used as a subroutine and its learning rate is set using the function  $G$  defined in Section IV-B. Optimized values for the parameters are  $z = 2.618$  and  $a = 2.0979$ .

*C. Bounding the range of the  $y_t$ 's*

We now introduce an algorithm **G1** for the case where a bound  $X$  on the norm of  $\mathbf{x}_t$  and a bound  $Y$  on  $|y_t|$ , for all  $t$ , are known ahead of time. The algorithm is sketched in Figure 2. In the following theorem we show a bound on the difference between the total loss of **G1** and the loss of the best linear predictor  $\mathbf{w}$  whose norm is bounded by  $Y/X$ , where  $X$  bounds the norm of the  $\mathbf{x}_t$ 's and  $Y$  bounds the norm of the  $y_t$ 's. The bound  $Y/X$  on the norm of the best linear predictor comes from an application of Theorem IV.3 and is the largest value for which we can prove the result.

*Theorem IV.4:* For each  $X, Y \geq 0$ , the algorithm **G1** has the following properties. Choose  $m \in \mathbf{N}$ ,  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_{t \leq m} \in (\mathcal{X} \times [-Y, Y])^m$  such that  $\max_t \|\mathbf{x}_t\| \leq X$ . Let  $\hat{y}_1, \dots, \hat{y}_m$  be the sequence of  $\mathbf{G1}_{X, Y}$ 's on-line predictions for  $\mathbf{s}$ . Then

$$\sum_{t=1}^m (\hat{y}_t - y_t)^2 \leq L_{Y/X}(\mathbf{s}) + 9.2 \left( Y \sqrt{L_{Y/X}(\mathbf{s})} + Y^2 \right).$$

**Proof of Theorem IV.4.** In Appendix A.

Our next result is a corollary to Theorem IV.4 for the normalized loss. We introduce a new algorithm, **G1-norm**, that differs from the algorithm **G1** only in the setting of the learning rate for the subroutine  $\mathbf{GD}$  (cfr. Figure 2.) That is, in each trial, **G1-norm** sets  $\mathbf{GD}$ 's learning rate to  $G(k_i, Y, 1)/\|\mathbf{x}_t\|^2$ . Thus, there is no need to know a bound  $X$  on the norm of the  $\mathbf{x}_t$ 's.

*Theorem IV.5:* For all  $Y \geq 0$ , the algorithm **G1-norm** has the following properties. Choose  $m \in \mathbf{N}$ ,  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_{t \leq m} \in (\mathcal{X} \times [-Y, Y])^m$ . Let  $\hat{y}_1, \dots, \hat{y}_m$  be the sequence of  $\mathbf{G1-norm}_Y$ 's on-line predictions for  $\mathbf{s}$ . Then

$$\sum_{t=1}^m \frac{(\hat{y}_t - y_t)^2}{\|\mathbf{x}_t\|^2} \leq L'_Y(\mathbf{s}) + 9.2 \left( Y \sqrt{L'_Y(\mathbf{s})} + Y^2 \right).$$

**Proof.** Given Corollary IV.1, the proof follows from a straightforward adaptation to the normalized loss of the proof of Theorem IV.4.  $\square$

*D. Predicting with no a priori information*

In this section we remove all assumptions that the learner has prior knowledge. We introduce a new variant of the  $\mathbf{GD}$  algorithm which we call **G2**. This new variant is described in Figure 3. A bound on **G2**'s total loss follows quite straightforwardly from Theorem IV.1 via the application of standard doubling techniques.

*Theorem IV.6:* For any  $0 < \beta < 2$ , the algorithm  $\mathbf{G2}_\beta$  has the following properties. Choose  $m \in \mathbf{N}$ ,  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_{t \leq m} \in (\mathcal{X} \times \mathbf{R})^m$ . Let  $\hat{y}_1, \dots, \hat{y}_m$  be the sequence of  $\mathbf{G2}_\beta$ 's on-line predictions for  $\mathbf{s}$ . Then,

$$\begin{aligned} \sum_{t=1}^m (\hat{y}_t - y_t)^2 &\leq \inf_{\mathbf{w} \in \mathcal{X}} \left[ \frac{4(\max_t \|\mathbf{x}_t\|^2) \|\mathbf{w}\|^2}{\beta(2-\beta)c} + \frac{L_{\mathbf{w}}(\mathbf{s})}{(2-\beta)^2 c(1-c)} \right] \end{aligned}$$

for all  $0 < c < 1$ . In particular, if  $\beta = 4/3$  and  $c = 1/2$ ,

$$\sum_{t=1}^m (\hat{y}_t - y_t)^2 \leq 9 \inf_{\mathbf{w} \in \mathcal{X}} \left[ (\max_t \|\mathbf{x}_t\|^2) \|\mathbf{w}\|^2 + L_{\mathbf{w}}(\mathbf{s}) \right].$$

**Proof:** Choose  $0 < \beta < 2$  and  $0 < c < 1$ . Notice that, in addition to a vector of hypothesized weights, **G2** maintains an integer  $j$  between trials. Before learning takes place,  $j$  is set to 0. After **G2** receives  $\mathbf{x}_1$ , it sets  $X_1 = \|\mathbf{x}_1\|$  and starts as a subroutine  $\mathbf{GD}_{\beta/(X_1)^2}$ . Thereafter, at each trial  $t$ , after **G2** receives  $\mathbf{x}_t$ , it sets

$$j \leftarrow \max \left\{ j, \left\lceil \log_{\sqrt{2}} \frac{\|\mathbf{x}_t\|}{X_1} \right\rceil \right\}.$$

Then **G2** uses  $\mathbf{GD}_{\beta/(2^{j/2} X_1)^2}$  for prediction on that trial.

Thus **G2** uses  $\mathbf{GD}_{\beta/(X_1)^2}$  as long as the  $\mathbf{x}_t$ 's are smaller than  $X_1$ , at which time it switches over to  $\mathbf{GD}_{\beta/(\sqrt{2} X_1)^2}$ , which it uses as long as the  $\mathbf{x}_t$ 's are no bigger than  $\sqrt{2} X_1$  (possibly for 0 trials), and continues in this manner, successively multiplying its assumed upper bound on the 2-norm of the  $\mathbf{x}_t$ 's by  $\sqrt{2}$ . Let  $X = \max_t \|\mathbf{x}_t\|$ . It follows immediately from Theorem IV.1 that

$$\begin{aligned} \sum_{t=1}^m (\hat{y}_t - y_t)^2 &\leq \left( \sum_{i=0}^{\lceil \log_{\sqrt{2}}(X/X_1) \rceil} \frac{\|\mathbf{w}\|^2 (2^{i/2} X_1)^2}{\beta(2-\beta)c} \right) + \frac{L_{\mathbf{w}}(\mathbf{s})}{(2-\beta)^2 c(1-c)} \\ &= \frac{\|\mathbf{w}\|^2 X_1^2}{\beta(2-\beta)c} \left( \sum_{i=0}^{\lceil \log_{\sqrt{2}}(X/X_1) \rceil} 2^i \right) + \frac{L_{\mathbf{w}}(\mathbf{s})}{(2-\beta)^2 c(1-c)} \\ &< \frac{\|\mathbf{w}\|^2 X_1^2}{\beta(2-\beta)c} 2^{2+2 \log_2(X/X_1)} + \frac{L_{\mathbf{w}}(\mathbf{s})}{(2-\beta)^2 c(1-c)} \\ &= \frac{4\|\mathbf{w}\|^2 X^2}{\beta(2-\beta)c} + \frac{L_{\mathbf{w}}(\mathbf{s})}{(2-\beta)^2 c(1-c)}. \end{aligned}$$

Plugging in  $\beta = 4/3$  and  $c = 1/2$  completes the proof.  $\square$

## V. APPLICATION TO CLASSES OF SMOOTH FUNCTIONS

In this section, we describe applications of the inner product results of the previous section to classes of smooth functions. While we will focus on applications of Theorem IV.3, we note that analogs of the other results of Section IV can be obtained in a similar manner.



**Algorithm G2.****Input**  $0 < \beta < 2$ .

- Let  $j = 0$ .
- Let  $X_1 = \|\mathbf{x}_1\|$ .
- On each trial  $t$ :
  1. Let  $j \leftarrow \max \left\{ j, \left\lceil \log_{\sqrt{2}} \frac{\|\mathbf{x}_t\|}{X_1} \right\rceil \right\}$ .
  2. Give  $\mathbf{x}_t$  to  $\mathbf{GD}_{\beta/(2^{j/2}X_1)^2}$ .
  3. Use  $\mathbf{GD}_{\beta/(2^{j/2}X_1)^2}$ 's prediction  $\hat{y}_t$ .
  4. Pass  $y_t$  to  $\mathbf{GD}_{\beta/(2^{j/2}X_1)^2}$ .

Fig. 3. Pseudo-code for the algorithm **G2** that uses **GD** as a subroutine. (See Theorem IV.6.) The learning rate of **GD** is dynamically set depending on the relative sizes of the  $\mathbf{x}_t$ 's.

*A. Smooth functions of a single variable*

We begin with a class of smooth functions of a single real variable that was studied by Faber and Mycielski [6] in a similar context, except using the assumption that there was a function  $f$  in the class such that  $y_t = f(x_t)$  for all  $t$ . Their methodology was to prove general results like those of the previous section under that assumption that there was a  $w$  with  $f_w(\mathbf{x}_t) = y_t$  for all  $t$ , then to reduce the smooth function learning problem to the more general problem as we do below. Similar function classes have also often been studied in nonparametric statistics (see, e.g. [9]) using probabilistic assumptions on the generation of the  $x_t$ 's.

Let  $\mathbf{R}_+$  be the set of nonnegative reals. We define the set  $\text{SMO}_W$  to be all absolutely continuous  $f : \mathbf{R}_+ \rightarrow \mathbf{R}$  for which

1.  $f(0) = 0$
2.  $\sqrt{\int_0^\infty f'(z)^2 dz} \leq W$ .

The assumption that  $f(0) = 0$  will be satisfied by many natural functions of interest. Examples include distance traveled as a function of time and return as a function of investment. We will prove the following result about  $\text{SMO}_W$ .

*Theorem V.1:* For each  $E, X, W \geq 0$ , there is a prediction algorithm  $\mathbf{A}_{\text{SMO}}$  with the following properties. Choose  $m \in \mathbf{N}$ ,  $\mathbf{s} = \langle (x_t, y_t) \rangle_{t \leq m} \in ([0, X] \times \mathbf{R})^m$ , such that there is an  $f \in \text{SMO}_W$  for which  $\sum_{t=1}^m (f(x_t) - y_t)^2 \leq E$ . Let  $\hat{y}_1, \dots, \hat{y}_m$  be the sequence of  $\mathbf{A}_{\text{SMO}}$ 's on-line predictions for  $\mathbf{s}$ . Then,

$$\sum_{t=1}^m (\hat{y}_t - y_t)^2 \leq \inf_{f \in \text{SMO}_W} \left[ \sum_{t=1}^m (f(x_t) - y_t)^2 \right] + 2W\sqrt{XE} + W^2X.$$

**Proof:** For now, let us ignore computational issues. We'll treat them again after the proof.

Fix  $E, X, W \geq 0$ . The algorithm  $\mathbf{A}_{\text{SMO}}$  operates by reducing the problem of learning  $\text{SMO}_W$  to a more general problem of the type treated in the previous section.

Let  $L^2(\mathbf{R}_+)$  be the space of (measurable) functions  $g$  from  $\mathbf{R}_+$  to  $\mathbf{R}$  for which  $\int_0^\infty g(u)^2 du$  is finite.  $L^2(\mathbf{R}_+)$  is well known to be an inner product space (see, e.g. [30]), with the inner product defined by

$$(g_1, g_2) = \int_0^\infty g_1(u)g_2(u) du.$$

**Algorithm A<sub>SMO</sub>.****Input:**  $E, W, X \geq 0$ .

- On each trial  $t$ :
  1. Get  $x_t \in [0, X]$  from the environment.
  2. Give  $\chi_{\leq x_t} \in L^2(\mathbf{R}_+)$  to  $\mathbf{GD}_{G(E,W,X)/X^2}$ .
  3. Use  $\mathbf{GD}_{G(E,W,X)/X^2}$ 's prediction  $\hat{y}_t$ .
  4. Pass  $y_t$  to  $\mathbf{GD}_{G(E,W,X)/X^2}$ .

Fig. 4. Pseudo-code for algorithm  $\mathbf{A}_{\text{SMO}}$ . (See Theorem V.1.) Algorithm **GD** (here used as a subroutine) is applied to the inner product space  $\mathcal{X} = L^2(\mathbf{R}_+)$ . The function  $G$ , used to set **GD**'s learning rate, is defined in Section IV-B.

Further, we define  $g_3 = g_2 + g_1$  by

$$(\forall x) \quad g_3(x) = g_2(x) + g_1(x),$$

and  $g_3 = \kappa g_1$  by

$$(\forall x) \quad g_3(x) = \kappa g_1(x).$$

Now apply algorithm **GD** to this particular inner product space,  $L^2(\mathbf{R}_+)$ , with learning rate  $\eta$  set to  $G(E, W, X)$ , where the function  $G$  is defined in Section IV-B. For any  $x \geq 0$ , define  $\chi_{\leq x} : \mathbf{R}_+ \rightarrow \mathbf{R}$  by

$$\chi_{\leq x}(u) = \begin{cases} 1 & \text{if } u \leq x \\ 0 & \text{otherwise.} \end{cases}$$

Note that for any  $x \leq X$

$$\|\chi_{\leq x}\| = \sqrt{\int_0^\infty \chi_{\leq x}(u)^2 du} = \sqrt{x} \leq \sqrt{X}, \quad (16)$$

and therefore  $\chi_{\leq x} \in L^2(\mathbf{R}_+)$ .

In Figure 4, we give a short description of the algorithm  $\mathbf{A}_{\text{SMO}}$ . Note that for any  $f \in \text{SMO}_W$ ,

$$\|f'\| = \sqrt{\int_0^\infty f'(u)^2 du} \leq W. \quad (17)$$

Finally, note that since  $f(0) = 0$ ,

$$(f', \chi_{\leq x}) = \int_0^\infty f'(u)\chi_{\leq x}(u) du = \int_0^x f'(u) du = f(x). \quad (18)$$

Thus, if there is an  $f \in \text{SMO}_W$  for which  $\sum_{t=1}^m (f(x_t) - y_t)^2 \leq E$ , then  $f' \in L^2(\mathbf{R}_+)$  has  $\|f'\| \leq W$  and satisfies

$$\sum_{t=1}^m ((f', \chi_{\leq x_t}) - y_t)^2 \leq E.$$

Combining this with (16) and Theorem IV.3, we can see that **GD**'s predictions satisfy

$$\sum_{t=1}^m (\hat{y}_t - y_t)^2 \leq \inf_{\|f'\| \leq W} \left[ \sum_{t=1}^m ((f', \chi_{\leq x_t}) - y_t)^2 \right] + 2W\sqrt{XE} + W^2X.$$

The result then follows from the fact that  $\mathbf{A}_{\text{SMO}}$  just makes the same predictions as **GD**.  $\square$

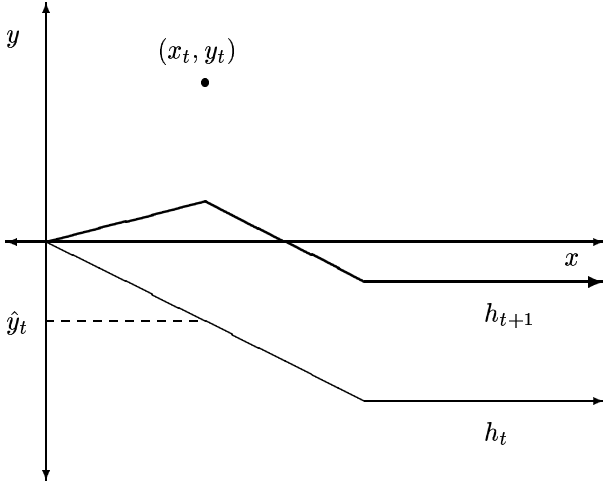


Fig. 5. An example of the update of the application of the **GD** algorithm to smoothing in the single-variable case. The derivative of the hypothesis is modified by a constant in the appropriate direction to the left of  $x_t$ , and left unchanged to the right.

By closely examining the predictions of the algorithm  $\mathbf{A}_{\text{SMO}}$  of Theorem V.1, we can see that it can be implemented in time polynomial in  $t$ . The algorithm **GD** maintains a function  $\hat{\mathbf{w}} \in L^2(\mathbf{R}_+)$  which it updates between trials. As before, let  $\hat{\mathbf{w}}_t$  be the  $t$ th hypothesis of **GD**. We can see that  $\hat{\mathbf{w}}_t$  can be interpreted as the derivative of  $\mathbf{A}_{\text{SMO}}$ 's  $t$ th hypothesis. This is because **GD**'s  $t$ th prediction, and therefore  $\mathbf{A}_{\text{SMO}}$ 's  $t$ th prediction, is

$$(\hat{\mathbf{w}}_t, \chi_{\leq x_t}) = \int_0^\infty \hat{\mathbf{w}}_t(u) \chi_{\leq x_t}(u) du = \int_0^{x_t} \hat{\mathbf{w}}_t(u) du.$$

Hence  $\mathbf{A}_{\text{SMO}}$ 's  $t$ th hypothesis  $h_t$  satisfies  $h_t' = \hat{\mathbf{w}}_t$ .

**GD** sets  $\hat{\mathbf{w}}_1$  to be the constant 0 function, and its update is

$$\hat{\mathbf{w}}_{t+1} = \hat{\mathbf{w}}_t + \eta(y_t - \hat{y}_t) \chi_{\leq x_t},$$

where  $\eta$  doesn't depend on  $t$  (see the proof of Theorem IV.3). Integrating yields the following expression for  $\mathbf{A}_{\text{SMO}}$ 's  $t+1$ st hypothesis:

$$h_{t+1}(x) = \begin{cases} h_t(x) + \eta(y_t - \hat{y}_t)x & \text{if } x \leq x_t \\ h_t(x) + \eta(y_t - \hat{y}_t)x_t & \text{otherwise} \end{cases}$$

and therefore

$$h_{t+1}(x) = h_t(x) + \eta(y_t - \hat{y}_t) \min\{x_t, x\}.$$

By induction, we have

$$h_{t+1}(x) = \eta \sum_{s \leq t} (y_s - \hat{y}_s) \min\{x_s, x\},$$

trivially computable in  $O(t)$  time if the previous  $\hat{y}_s$ 's are saved. This algorithm is illustrated in Figure 5.

### B. Smooth functions of several variables

Theorem V.1 can be generalized to higher dimensions as follows. The analogous generalization in the absence of

noise was carried out in [6]. The domain  $X$  is  $\mathbf{R}_+^n$ . We define the set  $\text{SMO}_{W,n}$  to be all functions  $f : \mathbf{R}_+^n \rightarrow \mathbf{R}$  for which there is a function  $\tilde{f}$  such that

1.  $\forall \mathbf{x} \in \mathbf{R}_+^n f(\mathbf{x}) = \int_0^{x_1} \dots \int_0^{x_n} \tilde{f}(u_1, \dots, u_n) du_n \dots du_1$
2.  $\sqrt{\int_0^\infty \dots \int_0^\infty (\tilde{f}(u_1, \dots, u_n))^2 du_n \dots du_1} \leq W$ .

It is easily verified that when  $\tilde{f}$  exists, it is defined by

$$\tilde{f}(u_1, \dots, u_n) = \frac{\partial^n f(u_1, \dots, u_n)}{\partial u_1 \dots \partial u_n}.$$

We can establish the following generalization of Theorem V.1.

*Theorem V.2:* For each  $E, X, W \geq 0$  and  $n \in \mathbf{N}$ , there is a prediction algorithm  $\mathbf{A}_{\text{SMO}n}$  with the following properties. Choose  $m \in \mathbf{N}$ ,  $\mathbf{s} = \langle (x_t, y_t) \rangle_{t \leq m} \in ([0, X]^n \times \mathbf{R})^m$ , such that there is an  $f \in \text{SMO}_{W,n}$  for which  $\sum_{t=1}^m (f(\mathbf{x}_t) - y_t)^2 \leq E$ . Let  $\hat{y}_1, \dots, \hat{y}_m$  be the sequence of  $\mathbf{A}_{\text{SMO}n}$ 's online predictions for  $\mathbf{s}$ . Then  $\sum_{t=1}^m (\hat{y}_t - y_t)^2$  is at most

$$\inf_{f \in \text{SMO}_{W,n}} \left[ \sum_{t=1}^m (f(\mathbf{x}_t) - y_t)^2 \right] + 2WX^{n/2}\sqrt{E} + W^2X^n.$$

**Proof.** Fix  $E, X, W, n \geq 0$ . The algorithm  $\mathbf{A}_{\text{SMO}n}$  operates by reducing the problem of learning  $\text{SMO}_{W,n}$  to a more general problem of the type treated in the previous section.

Let  $L^2(\mathbf{R}_+^n)$  be the space of (measurable) functions  $g$  from  $\mathbf{R}_+^n$  to  $\mathbf{R}$  for which

$$\int_0^\infty \dots \int_0^\infty g(\mathbf{x})^2 dx_n \dots dx_1$$

is finite. Again, it is well known (see e.g. [30]), that  $L^2(\mathbf{R}_+^n)$  has an inner product defined by

$$(g_1, g_2) = \int_0^\infty \dots \int_0^\infty g_1(\mathbf{x})g_2(\mathbf{x}) dx_n \dots dx_1.$$

Now apply algorithm **GD** to this particular inner product space,  $L^2(\mathbf{R}_+^n)$ , with learning rate  $\eta$  set to  $G(E, W, X)$ , where the function  $G$  is defined in Section IV-B. For any  $\mathbf{x} \in \mathbf{R}_+^n$ , define  $\chi_{\leq \mathbf{x}} : \mathbf{R}_+^n \rightarrow \mathbf{R}$  as the indicator function of the rectangle  $[0, x_1] \times \dots \times [0, x_n]$ . Note that for any  $\mathbf{x} \in [0, X]^n$

$$\|\chi_{\leq \mathbf{x}}\| = \sqrt{\prod_{i=1}^n x_i} \leq X^{n/2}. \quad (19)$$

and therefore  $\chi_{\leq \mathbf{x}} \in L^2(\mathbf{R}_+^n)$ .

The algorithm  $\mathbf{A}_{\text{SMO}n}$  is sketched in Figure 6. Note that for any  $f \in \text{SMO}_{W,n}$ , there is a function  $\tilde{f}$  such that

$$\begin{aligned} (\tilde{f}, \chi_{\leq \mathbf{x}_t}) &= \int_0^\infty \dots \int_0^\infty \tilde{f}(x_1, \dots, x_n) \chi_{\leq \mathbf{x}_t}(x_1, \dots, x_n) dx_n \dots dx_1 \\ &= f(\mathbf{x}_t). \end{aligned}$$

Thus, if there is an  $f \in \text{SMO}_{W,n}$  for which  $\sum_{t=1}^m (f(\mathbf{x}_t) - y_t)^2 \leq E$ , then the corresponding  $\tilde{f} \in L^2(\mathbf{R}_+^n)$ , which has  $\|\tilde{f}\| \leq W$ , satisfies  $\sum_{t=1}^m ((\tilde{f}, \chi_{\leq \mathbf{x}_t}) - y_t)^2 \leq E$ . Combining

**Algorithm  $\mathbf{A}_{\text{SMOn}}$ .****Input:**  $E, W, X \geq 0$ .

- On each trial  $t$ :
  1. Get  $\mathbf{x}_t \in [0, X]^n$  from the environment.
  2. Give  $\chi_{\leq \mathbf{x}_t} \in L^2(\mathbf{R}_+^n)$  to  $\mathbf{GD}_{G(E,W,X)/X^2}$ .
  3. Use  $\mathbf{GD}_{G(E,W,X)/X^2}$ 's prediction  $\hat{y}_t$ .
  4. Pass  $y_t$  to  $\mathbf{GD}_{G(E,W,X)/X^2}$ .

Fig. 6. Pseudo-code for algorithm  $\mathbf{A}_{\text{SMOn}}$ . (See Theorem V.2.) Algorithm  $\mathbf{GD}$  (here used as a subroutine) is applied to the inner product space  $\mathcal{X} = L^2(\mathbf{R}_+^n)$ . The function  $G$ , used to set  $\mathbf{GD}$ 's learning rate, is defined in Section IV-B.

this with (19) and Theorem IV.3, we can see that, for  $\mathbf{GD}$ 's predictions,  $\sum_{t=1}^m (\hat{y}_t - y_t)^2$  is bounded by

$$\inf_{\|\tilde{f}\| \leq W} \left[ \sum_{t=1}^m ((\tilde{f}, \chi_{\leq \mathbf{x}_t}) - y_t)^2 \right] + 2WX^{n/2}\sqrt{E} + W^2X^n.$$

The result then follows from the fact that  $\mathbf{A}_{\text{SMOn}}$  just makes the same predictions as  $\mathbf{GD}$ .  $\square$

It is easy to see, by extending the discussion following Theorem V.1, that the predictions of Theorem V.2 can be computed in  $O(tn)$  time, if previous predictions are saved.

## VI. A COMPARISON TO STANDARD GRADIENT DESCENT METHODS

The goal of this section is to compare the total square loss bounds obtained via our analysis to the bounds obtained via the standard analysis of gradient descent methods. Standard methods only deal with the case when all the pairs  $(\mathbf{x}_t, y_t)$  are given at once (batch case) rather than in an on-line fashion. Thus we consider the problem of finding the solution  $\mathbf{x} \in \mathbf{R}^n$  of a system of linear equations

$$\begin{aligned} a_{1,1}x_1 + a_{1,2}x_2 + \dots + a_{1,n}x_n &= b_1 \\ &\vdots \\ a_{m,1}x_1 + a_{m,2}x_2 + \dots + a_{m,n}x_n &= b_m \end{aligned}$$

where  $a_{i,j}, b_i \in \mathbf{R}$ . The above system can be given the more compact representation  $A\mathbf{x} = \mathbf{b}$ , where  $\mathbf{b} = (b_1, \dots, b_m)$  and  $A$  is a  $m \times n$  matrix with entries  $a_{i,j}$ . ( $A\mathbf{x}$  denotes the usual matrix-vector product.) For simplicity, we assume in this section that  $A\mathbf{x} = \mathbf{b}$  has a solution. However, we do not assume that the matrix  $A$  has any special property.

A standard iterative approach for solving the problem  $A\mathbf{x} = \mathbf{b}$  is to perform gradient descent over the (total) squared residual error  $R(\mathbf{x}) = \|A\hat{\mathbf{x}} - \mathbf{b}\|_2^2$ , where  $\hat{\mathbf{x}}$  is a candidate solution. We will prove upper bounds on the sum of  $R(\hat{\mathbf{x}}_t)$  for the sequence  $\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, \dots$  of candidate solutions generated by the gradient descent method tuned either according to the standard analysis or to our analysis. The bounds are expressed in terms of both the norm of the solution  $\mathbf{x}$  and the eigenvalues of  $A^T A$ , where  $A^T$  denotes the transpose matrix of  $A$ .

We define the norm  $\|A\|$  of a matrix  $A$  by

$$\|A\|_2 = \sup_{\|\mathbf{v}\|_2=1} \|A\mathbf{v}\|_2.$$

This is the norm induced by the Euclidean norm for vectors in  $\mathbf{R}^n$  (see [8].) Notice that  $\|A\mathbf{v}\|_2 \leq \|A\|_2 \|\mathbf{v}\|_2$  (Cauchy-Schwartz inequality). We will make use of the following well-known facts.

*Fact VI.1* ([11]) For any real matrix  $A$ ,  $\|A\|_2 = \sqrt{\lambda_{\max}}$ , where  $\lambda_{\max}$  is the largest eigenvalue of  $A^T A$ .

*Fact VI.2* ([11]) For any real matrix  $A$ ,

$$\|A^T\|_2 = \|A\|_2.$$

Given a candidate solution  $\hat{\mathbf{x}} \in \mathbf{R}^n$  with squared residual error  $R(\hat{\mathbf{x}})$ , the gradient of  $R(\hat{\mathbf{x}})$  with respect to  $\hat{\mathbf{x}}$  is  $\vec{\nabla} R(\hat{\mathbf{x}}) = 2A^T(A\hat{\mathbf{x}} - \mathbf{b})$ . By applying the gradient descent (Kaczmarz) rule for the batch case we derive the update

$$\hat{\mathbf{x}}_{t+1} = \hat{\mathbf{x}}_t - \eta 2A^T(A\hat{\mathbf{x}}_t - \mathbf{b}) \quad (20)$$

for some scaling factor  $\eta > 0$ . Simple manipulation shows that

$$R(\hat{\mathbf{x}}_{t+1}) = R(\hat{\mathbf{x}}_t) + \eta^2 \|A\vec{\nabla} R(\hat{\mathbf{x}}_t)\|_2^2 - \eta \|\vec{\nabla} R(\hat{\mathbf{x}}_t)\|_2^2. \quad (21)$$

Following the standard analysis of gradient descent, we find the value of  $\eta$  minimizing the LHS of (21) as

$$\eta_1 = \frac{\|\vec{\nabla} R(\hat{\mathbf{x}}_t)\|_2^2}{2\|A\vec{\nabla} R(\hat{\mathbf{x}}_t)\|_2^2}.$$

By plugging this optimal value of  $\eta$  back in (21) we get

$$R(\hat{\mathbf{x}}_{t+1}) = R(\hat{\mathbf{x}}_t) - \frac{\|\vec{\nabla} R(\hat{\mathbf{x}}_t)\|_2^4}{4\|A\vec{\nabla} R(\hat{\mathbf{x}}_t)\|_2^2}.$$

*Proposition VI.1:* For all  $m, n > 0$ , for any  $m \times n$  real matrix  $A$  and for any vector  $\mathbf{x} \in \mathbf{R}^n$ . Let  $\mathbf{b} = A\mathbf{x}$  and let  $\lambda_{\min}, \lambda_{\max}$  be, respectively, the smallest and the largest eigenvalues of  $A^T A$ . Then, if  $\hat{\mathbf{x}}_0 = \mathbf{0}$  and  $\hat{\mathbf{x}}_{t+1}$  is computed from  $\hat{\mathbf{x}}_t$  using formula (20) with  $\eta = \eta_1$ ,

$$\sum_{t=0}^{\infty} \|A\hat{\mathbf{x}}_t - \mathbf{b}\|_2^2 \leq \frac{(\lambda_{\min} + \lambda_{\max})^2}{4\lambda_{\min}} \|\mathbf{x}\|_2^2.$$

**Proof.** If  $\lambda_{\min} = 0$ , then the bound holds vacuously. Assume then  $\lambda_{\min} > 0$ . Via an application of the Kantorovich inequality to the square matrix  $A^T A$  (see e.g. [20]) it can be shown that

$$R(\hat{\mathbf{x}}_{t+1}) \leq \left(1 - \frac{4\lambda_{\min}\lambda_{\max}}{(\lambda_{\min} + \lambda_{\max})^2}\right) R(\hat{\mathbf{x}}_t). \quad (22)$$

Therefore, we get

$$\frac{4\lambda_{\min}\lambda_{\max}}{(\lambda_{\min} + \lambda_{\max})^2} R(\hat{\mathbf{x}}_t) \leq R(\hat{\mathbf{x}}_t) - R(\hat{\mathbf{x}}_{t+1}).$$

By summing up over all iterations  $t$  we obtain

$$\frac{4\lambda_{\min}\lambda_{\max}}{(\lambda_{\min} + \lambda_{\max})^2} \sum_{t=0}^{\infty} R(\hat{\mathbf{x}}_t) \leq R(\hat{\mathbf{x}}_0).$$

Recalling that  $\hat{\mathbf{x}}_0 = (0, \dots, 0)$  and making use of Fact VI.1, □

$$\begin{aligned}
\sum_{t=0}^{\infty} \|A\hat{\mathbf{x}}_t - \mathbf{b}\|_2^2 &\leq \frac{(\lambda_{\min} + \lambda_{\max})^2}{4\lambda_{\min}\lambda_{\max}} R(\hat{\mathbf{x}}_0) \\
&\leq \frac{(\lambda_{\min} + \lambda_{\max})^2}{4\lambda_{\min}\lambda_{\max}} \|A\mathbf{x}\|_2^2 \\
&\leq \frac{(\lambda_{\min} + \lambda_{\max})^2}{4\lambda_{\min}\lambda_{\max}} \|A\|_2^2 \|\mathbf{x}\|_2^2 \\
&\leq \frac{(\lambda_{\min} + \lambda_{\max})^2}{4\lambda_{\min}\lambda_{\max}} \lambda_{\max} \|\mathbf{x}\|_2^2 \\
&= \frac{(\lambda_{\min} + \lambda_{\max})^2}{4\lambda_{\min}} \|\mathbf{x}\|_2^2
\end{aligned}$$

concluding the proof. □

A different analysis of update (20) can be obtained by applying the techniques developed in Section IV. Let  $D(\hat{\mathbf{x}})$  be the distance  $\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2$  of  $\hat{\mathbf{x}}$  to the solution  $\mathbf{x}$ . An easy adaptation of Lemma IV.1 shows that

$$D(\hat{\mathbf{x}}_{t+1}) = D(\hat{\mathbf{x}}_t) + \eta^2 \|\vec{\nabla} R(\hat{\mathbf{x}}_t)\|_2^2 - 4\eta R(\hat{\mathbf{x}}_t). \quad (23)$$

Here, the minimization over  $\eta$  yields the optimum at

$$\eta_2 = \frac{2R(\hat{\mathbf{x}}_t)}{\|\vec{\nabla} R(\hat{\mathbf{x}}_t)\|_2^2}.$$

We then have the following result.

*Proposition VI.2:* For all  $m, n > 0$ , for any  $m \times n$  real matrix  $A$  and for any vector  $\mathbf{x} \in \mathbf{R}^n$ . Let  $\mathbf{b} = A\mathbf{x}$  and let  $\lambda_{\max}$  be the largest eigenvalue of  $A^T A$ . Then, if  $\hat{\mathbf{x}}_0 = 0$  and  $\hat{\mathbf{x}}_{t+1}$  is computed from  $\hat{\mathbf{x}}_t$  using formula (20) with  $\eta = \eta_2$ ,

$$\sum_{t=0}^{\infty} \|A\hat{\mathbf{x}}_t - \mathbf{b}\|_2^2 \leq \lambda_{\max} \|\mathbf{x}\|_2^2.$$

**Proof.** By plugging  $\eta_2$  for  $\eta$  in (23) we obtain

$$\begin{aligned}
D(\hat{\mathbf{x}}_{t+1}) &= D(\hat{\mathbf{x}}_t) - \frac{4R(\hat{\mathbf{x}}_t)^2}{\|\vec{\nabla} R(\hat{\mathbf{x}}_t)\|_2^2} \\
&= D(\hat{\mathbf{x}}_t) - \|A\hat{\mathbf{x}}_t - \mathbf{b}\|_2^2 \frac{\|A\hat{\mathbf{x}}_t - \mathbf{b}\|_2^2}{\|A^T(A\hat{\mathbf{x}}_t - \mathbf{b})\|_2^2} \\
&\leq D(\hat{\mathbf{x}}_t) - \frac{\|A\hat{\mathbf{x}}_t - \mathbf{b}\|_2^2}{\|A^T\|_2} \\
&\quad (\text{by definition of } \|A^T\|_2) \\
&\leq D(\hat{\mathbf{x}}_t) - \frac{\|A\hat{\mathbf{x}}_t - \mathbf{b}\|_2^2}{\|A\|_2^2} \quad (\text{by Fact VI.2}).
\end{aligned}$$

Therefore, rearranging the above and summing up over all iterations  $t$ ,

$$\begin{aligned}
\sum_{t=0}^{\infty} \|A\hat{\mathbf{x}}_t - \mathbf{b}\|_2^2 &\leq \|A\|_2^2 D(\hat{\mathbf{x}}_0) \\
&= \|A\|_2^2 \|\mathbf{x}\|_2^2
\end{aligned}$$

since  $\hat{\mathbf{x}}_0 = (0, \dots, 0)$ . By Fact VI.1, this implies

$$\sum_{t=0}^{\infty} \|A\hat{\mathbf{x}}_t - \mathbf{b}\|_2^2 \leq \lambda_{\max} \|\mathbf{x}\|_2^2.$$

In summary, we compared two tunings of  $\eta$  for the learning rule (20). The first and standard one maximizes the decrease of  $\|A\hat{\mathbf{x}} - \mathbf{b}\|_2^2$  and the second one maximizes the decrease in  $\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2$ , where  $\mathbf{x}$  is a solution.

The first method has the advantage that one can show that  $\|A\hat{\mathbf{x}} - \mathbf{b}\|_2^2$  decreases by a fixed factor in each trial (Inequality (22)). (Note that this factor is 1 when  $\lambda_{\min} = 0$ , and this holds when  $A$  does not have full rank.) In contrast, matrices  $A$  can be constructed where updating with the optimal learning rate  $\eta_2$  causes an increase in  $\|A\hat{\mathbf{x}} - \mathbf{b}\|_2^2$ .

The second method, however, always leads to better bounds on  $\sum_t \|A\hat{\mathbf{x}}_t - \mathbf{b}\|_2^2$  since

$$\lambda_{\max} \leq \frac{(\lambda_{\min} + \lambda_{\max})^2}{4\lambda_{\min}}$$

for all  $\lambda_{\min}, \lambda_{\max} \geq 0$ . (Notice that the corresponding bound for the first method is vacuous when  $\lambda_{\min} = 0$ , which holds, as we said above, when  $A$  does not have full rank.)

## VII. LOWER BOUNDS

In this section, we describe lower bounds which match the upper bounds of Theorems IV.3, V.1, and V.2, constants included. In fact, these lower bounds show that even the upper bound on the excess of the algorithm's squared loss above the best fixed element within a given class of functions is optimal.

*Theorem VII.1:* Fix an inner product space  $\mathcal{X}$  for which an orthonormal basis can be found.<sup>6</sup> For all  $E, X, W \geq 0$  and all prediction algorithms  $\mathbf{A}$ , there exists  $n \in \mathbf{N}$  and a pair  $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathbf{R}$ , such that  $\|\mathbf{x}\| \leq X$  and the following hold: There is a  $\mathbf{w} \in \mathcal{X}$  for which  $\|\mathbf{w}\| = W$  and  $((\mathbf{w}, \mathbf{x}) - \mathbf{y})^2 = E$ . Furthermore, if  $\hat{\mathbf{y}} = A(\mathbf{x}_t)$  then

$$(\hat{\mathbf{y}} - \mathbf{y})^2 \geq E + 2WX\sqrt{E} + (WX)^2.$$

**Proof.** Choose an orthonormal basis for  $\mathcal{X}$ . Set  $\mathbf{x} = (X, 0, \dots)$ ,  $\mathbf{y} = \text{sgn}(-\hat{\mathbf{y}})(WX + \sqrt{E})$ , and  $\mathbf{w} = (\text{sgn}(-\hat{\mathbf{y}})W, 0, \dots)$ . The result then follows easily. □

To establish the upper bound of Theorem IV.4, in which general bounds were obtained without any knowledge of an upper bound on  $L_W(\mathbf{s})$ , we required the assumption that the  $\mathbf{y}_t$ 's were in a known range  $[-Y, Y]$  and compared the total loss of the **GD** algorithm on  $\mathbf{s}$  against  $L_W(\mathbf{s})$ , where  $W = Y/(\max_t \|\mathbf{x}_t\|)$ . Therefore, the above lower bound does not say anything about the optimality of those results. The following lower bound shows that Theorem IV.4 cannot be significantly improved, at least for high-dimensional spaces. That is, we show that for any given  $X, Y, E > 0$  there is some  $n_0$  such that for all  $n$ -dimensional inner product spaces  $\mathcal{X}_n$ , with  $n \geq n_0$ , any prediction strategy incurs loss at least  $E + 2Y\sqrt{E} + Y^2$  for some sequence on  $\mathcal{X}_n \times [-Y, Y]$ . This theorem further has obvious consequences concerning the finite dimension case

<sup>6</sup>An orthonormal basis can be found under quite general conditions. See e.g. [30] for details.

when the “noise level”  $E$  is not too large relative to the number  $n$  of variables as well as  $X$  and  $Y$ .

*Theorem VII.2:* Let  $\langle \mathcal{X}_d \rangle_{d \in \mathbf{N}}$  be any sequence of inner product spaces such that  $\mathcal{X}_d$  is a  $d$ -dimensional vector space. Choose  $X, Y, E > 0$ . Let  $n$  be any integer such that

$$n \geq \left(1 + \frac{\sqrt{E}}{Y}\right)^2. \quad (24)$$

Then for any prediction algorithm  $A$  there is a sequence  $\langle (\mathbf{x}_1, y_1) \rangle_{t \leq n} \in (\mathcal{X}_n \times [-Y, Y])^n$  such that

1. For all  $1 \leq t \leq n$ ,  $\|\mathbf{x}_t\| = X$ .
2. If for each  $t$ ,  $\hat{y}_t = A(\langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1}) \rangle, \mathbf{x}_t)$ , then

$$\sum_{t=1}^n (y_t - \hat{y}_t)^2 \geq (Y + \sqrt{E})^2 = E + 2Y\sqrt{E} + Y^2.$$

3. There exists  $\mathbf{w} \in \mathbf{R}^n$  such that  $\|\mathbf{w}\| = Y/X$  and

$$\sum_{t=1}^n (y_t - (\mathbf{w}, \mathbf{x}_t))^2 = E.$$

**Proof.** Choose  $X, Y, E > 0$  and choose  $n \in \mathbf{N}$  so that (24) is satisfied. Let  $\mathbf{e}_1, \dots, \mathbf{e}_n$  be an orthonormal basis of  $\mathcal{X}_n$  (since  $\mathcal{X}_n$  is a finite-dimensional inner product space, such an orthonormal basis can always be found). Let  $\mathbf{x}_i = X\mathbf{e}_i$ , for  $i = 1, \dots, n$ . Since the basis is orthonormal,  $\|\mathbf{x}_i\| = X$  for all  $i$ , fulfilling part 1. Consider the adversary which at each step  $t = 1, \dots, n$  feeds the algorithm with vector  $\mathbf{x}_t$  and, upon algorithm’s prediction  $\hat{y}_t$ , responds with

$$y_t = \text{sgn}(-\hat{y}_t) \frac{Y + \sqrt{E}}{\sqrt{n}}.$$

This implies

$$(y_t - \hat{y}_t)^2 \geq \left(\frac{Y + \sqrt{E}}{\sqrt{n}}\right)^2$$

for all  $t = 1, 2, \dots, n$ . This proves part 2. Now let  $\mathbf{w}$  be the vector of  $\mathcal{X}_n$  with coordinates

$$\left(\text{sgn}(-\hat{y}_1) \frac{Y/X}{\sqrt{n}}, \dots, \text{sgn}(-\hat{y}_n) \frac{Y/X}{\sqrt{n}}\right)$$

with respect to the basis  $\mathbf{e}_1, \dots, \mathbf{e}_n$ . To prove part 3, first notice that  $\|\mathbf{w}\| = Y/X$ . Second, for each  $t = 1, \dots, n$  we have

$$\begin{aligned} & (y_t - (\mathbf{x}_t, \mathbf{w}))^2 \\ &= \left[ \text{sgn}(-\hat{y}_t) \frac{Y + \sqrt{E}}{\sqrt{n}} - (\mathbf{x}_t, \mathbf{w}) \right]^2 \\ &= \left[ \text{sgn}(-\hat{y}_t) \frac{Y + \sqrt{E}}{\sqrt{n}} - X(\mathbf{e}_t, \mathbf{w}) \right]^2 \\ &= \left[ \text{sgn}(-\hat{y}_t) \frac{Y + \sqrt{E}}{\sqrt{n}} - X \left( \text{sgn}(-\hat{y}_t) \frac{Y/X}{\sqrt{n}} \right) \right]^2 \\ &= \left( \frac{Y + \sqrt{E}}{\sqrt{n}} - \frac{Y}{\sqrt{n}} \right)^2 \\ &= \frac{E}{n}. \end{aligned}$$

This concludes the proof of part 3. Finally, notice that (24) implies that for all  $t = 1, 2, \dots, n$

$$|y_t| = \frac{Y + \sqrt{E}}{\sqrt{n}} \leq Y.$$

The proof is complete.  $\square$

We conclude with a lower bound for smooth functions.

*Theorem VII.3:* Choose  $E, X, W \geq 0$ ,  $n \in \mathbf{N}$ , and a prediction algorithm  $\mathbf{A}$ . Then there exists  $m \in \mathbf{N}$ ,  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_{t \leq m} \in ([0, X]^n \times \mathbf{R})^m$ , such that the following hold: There is a function  $f \in \text{SMO}_{W,n}$  for which  $\sum_{t=1}^m (f(\mathbf{x}_t) - y_t)^2 \leq E$ . If for each  $t$ ,

$$\hat{y}_t = A(\langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1}) \rangle, \mathbf{x}_t),$$

then

$$\sum_{t=1}^m (\hat{y}_t - y_t)^2 \geq E + 2WX^{n/2}\sqrt{E} + 2W^2X^n$$

**Proof.** In fact  $m = 1$  suffices in this case. Let  $\mathbf{x}_1 = (X, \dots, X)$ . Suppose the first prediction  $\hat{y}_1$  of  $\mathbf{A}$  is non-positive. Let

$$y_1 = WX^{n/2} + \sqrt{E}$$

and let the function  $f: \mathbf{R}_+^n \rightarrow \mathbf{R}$  be defined by

$$f(\mathbf{x}) = \frac{W}{X^{n/2}} \prod_{i=1}^n x_i,$$

if  $\mathbf{x} \in [0, X]^n$ , and  $f(\mathbf{x}) = 0$  otherwise. Then, for any  $\mathbf{x} \in [0, X]^n$ ,

$$f(\mathbf{x}) = \int_0^{x_1} \dots \int_0^{x_n} \tilde{f}(u_1, \dots, u_n) du_n \dots du_1,$$

where  $\tilde{f} \equiv \frac{W}{X^{n/2}}$ . The following are then easily verified

1.  $f(0) = 0$
2.  $(f(\mathbf{x}_1) - y_1)^2 = (WX^{n/2} - (WX^{n/2} + \sqrt{E}))^2 = E$
3.  $\sqrt{\int_0^\infty \dots \int_0^\infty \tilde{f}(\mathbf{u})^2 du_n \dots du_1} = \sqrt{X^n (W/X^{n/2})^2} = W$
4.  $(\hat{y}_1 - y_1)^2 \geq (WX^{n/2} + \sqrt{E})^2 = E + 2WX^{n/2}\sqrt{E} + W^2X^n$

since  $\hat{y}_1 \leq 0$ . The case in which  $\hat{y}_1 > 0$  can be handled symmetrically.  $\square$

## VIII. DISCUSSION AND CONCLUSIONS

In this paper we have investigated the performance of the gradient descent rule applied to the problem of on-line prediction in arbitrary inner product spaces. Through a reduction, we then applied our results to natural classes of smooth functions.

One of the most interesting contributions of this work is perhaps the derivation of the optimal “learning rate” for gradient descent methods when the goal is to minimize the worst-case total loss (here the sum of the squared prediction errors). Our tuning of the learning rate is based on *a priori* information that can be guessed on-line with an increase in the total loss of constant factors only. In the case

of iterative solution of systems of linear equations, we also showed that, with respect to the sum of squared residual errors, the tuning provided by our analysis compares favorably against the tuning obtained via the standard gradient descent analysis.

It is an open problem whether, instead of using adversarial arguments as we do here, our lower bounds can already be obtained when the examples are randomly and independently drawn from a natural distribution. For more simple functions this was done in [3]: the lower bounds there are with respect to uniform distributions and the upper bounds which essentially meet the lower bounds are proven for the worst-case as done in this paper.

An interesting open problem is whether a variant of the  $\mathbf{GD}_{X,Y}$  algorithm (see Figure 2) exists such that, for all sequences  $\mathbf{s} = \langle (\mathbf{x}_t, y_t) \rangle_{t \leq m}$  satisfying  $\|\mathbf{x}_t\| \leq X$  and  $|y_t| \leq Y$  for all  $t$ , the additional total loss of the algorithm on  $\mathbf{s}$  over and above  $\inf_{\mathbf{w} \in \mathcal{X}} L\mathbf{w}(\mathbf{s})$  is bounded by a function of  $X, Y$  only. Notice that this does not contradict Theorem VII.2.

The most challenging research direction is to prove worst case loss bounds for other gradient descent applications (by tuning the learning rate) as we have done in this paper for linear functions and the square loss. For example, are there useful worst case loss bounds for learning linear functions with other loss functions than the square loss. Another interesting case would be worst case loss bounds for learning the class of linear functions passed through a fixed transfer function (such as tanh or the sigmoid function) for any reasonable loss function.

### Acknowledgements

We thank Ethan Bernstein for helpful conversations, and for pointing out an error in an earlier version of this paper. Thanks to Jyrki Kivinen for simplifying the proof of Theorem VII.2 and to an anonymous referee for simplifying the proof of Lemma 8. Thanks also to Peter Auer, Gianpiero Cattaneo, and Shuxian Lou for their comments. Finally, we are grateful to Jan Mycielski for telling us about the paper by Kaczmarsz.

### REFERENCES

- [1] D. Angluin. Queries and concept learning. *Machine Learning*, 2:319–342, 1988.
- [2] A. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transactions on Information Theory*, 39(3):930–945, 1993.
- [3] N. Cesa-Bianchi, Y. Freund, D.P. Helmbold, D. Haussler, R.E. Schapire, and M.K. Warmuth. How to use expert advice. *Proceedings of the 25th ACM Symposium on the Theory of Computation*, 1993.
- [4] A. Dawid. Statistical theory: The prequential approach. *Journal of the Royal Statistical Society (Series A)*, pages 278–292, 1984.
- [5] R. O. Duda and P. E. Hart. *Pattern Classification and Scene Analysis*. Wiley, 1973.
- [6] V. Faber and J. Mycielski. Applications of learning theorems. *Fundamenta Informaticae*, 15(2):145–167, 1991.
- [7] M. Feder, N. Merhav, and M. Gutman. Universal prediction of individual sequences. *IEEE transactions of information theory*, 38:1258–1270, 1992.
- [8] G.H. Golub and C.F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, 1990.
- [9] W. Hardle. *Smoothing Techniques*. Springer Verlag, 1991.

- [10] S.S. Haykin. *Adaptive Filter Theory*. Prentice Hall, 2nd edition, 1991.
- [11] R.A. Horn and C.R. Johnson. *Matrix Analysis*. Cambridge University Press, 1985.
- [12] S. Hui and S.H. Zak. Robust stability analysis of adaptation algorithms for single perceptron. *IEEE Transactions on Neural Networks*, 2(2):325–328, 1991.
- [13] S. Kaczmarsz. Angenaherte Auflösung von systemen linearer gleichungen. *Bull. Acad. Polon. Sci. Lett. A*, 35:355–357, 1937.
- [14] D. Kimber and P.M. Long. The learning complexity of smooth functions of a single variable. *The 1992 Workshop on Computational Learning Theory*, pages 153–159, 1992.
- [15] J. Kivinen and M. K. Warmuth. Exponentiated gradient versus gradient descent for linear predictors. Report UCSC-CRL-94-16, University of California, Santa Cruz, June 1994. An extended abstract appeared in Proceedings of STOC 95.
- [16] N. Littlestone. Learning quickly when irrelevant attributes abound: a new linear-threshold algorithm. *Machine Learning*, 2(4):285–318, 1988.
- [17] N. Littlestone. *Mistake Bounds and Logarithmic Linear-threshold Learning Algorithms*. PhD thesis, Technical Report UCSC-CRL-89-11, University of California Santa Cruz, 1989.
- [18] N. Littlestone, P.M. Long, and M.K. Warmuth. On-line learning of linear functions. *Proceedings of the 23rd ACM Symposium on the Theory of Computation*, pages 465–475, 1991.
- [19] R.W. Lucky. Techniques of adaptive equalization of digital communication systems. *Bell System Technical Journal*, 45:255–286, 1966.
- [20] D.G. Luenberger. *Linear and Nonlinear Programming*. Addison-Wesley, 1984.
- [21] N. Littlestone and M.K. Warmuth. The weighted majority algorithm. Technical Report UCSC-CRL-91-28, UC Santa Cruz, October 1991. A preliminary version appeared in the *Proceedings of the 30th Annual IEEE Symposium on the Foundations of Computer Science*, October 89, pages 256–261.
- [22] N. Merhav and M. Feder. Universal sequential learning and decision from individual data sequences. *The 1992 Workshop on Computational Learning Theory*, pages 413–427, 1992.
- [23] J. Mycielski and S. Swierczkowski. General learning theorems. Unpublished, 1991.
- [24] J. Mycielski. A learning algorithm for linear operators. *Proceedings of the American Mathematical Society*, 103(2):547–550, 1988.
- [25] M.M. Sondhi. An adaptive echo canceller. *Bell System Technical Journal*, 46:497–511, 1967.
- [26] M.M. Sondhi and D.A. Berkley. Silencing echoes in the telephone network. *Proceedings of the IEEE*, 68:948–963, 1980.
- [27] V. Vovk. Aggregating strategies. In *Proceedings of the 3rd Workshop on Computational Learning Theory*, pages 371–383. Morgan Kaufmann, 1990.
- [28] B. Widrow and M.E. Hoff. Adaptive switching circuits. *1960 IRE WESCON Conv. Record*, pages 96–104, 1960.
- [29] B. Widrow and S.D. Stearns. *Adaptive signal processing*. Prentice Hall, 1985.
- [30] N. Young. *An introduction to Hilbert space*. Cambridge University Press, 1988.

### APPENDIX

#### I. PROOF OF THEOREM IV.4

Before proving the theorem we need some preliminary lemmas.

*Lemma A.1:* The total loss of  $\mathbf{G1}$  incurred in each loop  $i$  is at most  $k_i + (2az^{i/2} + 5)Y^2$ .

**Proof.** By construction of  $\mathbf{G1}$ , the total loss incurred in each loop  $i$  is at most  $k_i + (2az^{i/2} + 1)Y^2$  plus the possible additional loss on the trial causing the exit from the loop. To upper bound this additional loss observe that  $\mathbf{G1}$  always predicts with a value  $\hat{y}_t$  in the range  $[-Y, Y]$ . By hypothesis,  $y_t \in [-Y, Y]$  for all  $t$ . Hence the loss of  $\mathbf{G1}$  on a single trial  $t$  is at most  $4Y^2$ .  $\square$

In what follows  $W = Y/X$ . Let  $\mathbf{s}_i$  be the subsequence of  $\mathbf{s}$  fed to  $\mathbf{G1}$  during loop  $i$ .

*Lemma A.2:* If **G1** exits loop  $i$ , then  $L_W(\mathbf{s}_i) > k_i$ .

**Proof.** By construction of **G1**, if **G1** exits loop  $i$ , then the total loss incurred on subsequence  $\mathbf{s}_i$  is bigger than

$$k_i + 2Y\sqrt{k_i} + Y^2.$$

Since  $|y_t| \leq Y$  and since **G1** predicts on each trial of loop  $i$  by “clipping” the prediction of  $\mathbf{GD}_{G(k_i, W, X)/X^2}$  to make it fit in the range  $[-Y, Y]$ , we conclude that the total loss incurred by  $\mathbf{GD}_{G(k_i, W, X)/X^2}$  on loop  $i$  is bigger than  $k_i + 2Y\sqrt{k_i} + Y^2$  as well. Hence by Theorem IV.3  $L_W(\mathbf{s}_i) > k_i$  must hold.  $\square$

*Lemma A.3:* Let  $\ell$  be the index of the last loop entered by **G1**. Then

$$\ell \leq \log_z \left( 1 + \frac{(z-1)L_W(\mathbf{s})}{(aY)^2} \right).$$

**Proof.**

$$\begin{aligned} L_W(\mathbf{s}) &= \inf_{\|\mathbf{w}\| \leq W} L_{\mathbf{w}}(\mathbf{s}) \\ &= \inf_{\|\mathbf{w}\| \leq W} \left[ \sum_{i=0}^{\ell} L_{\mathbf{w}}(\mathbf{s}_i) \right] \\ &\geq \sum_{i=0}^{\ell} \left[ \inf_{\|\mathbf{w}\| \leq W} L_{\mathbf{w}}(\mathbf{s}_i) \right] \\ &= \sum_{i=0}^{\ell} L_W(\mathbf{s}_i) \\ &\geq \sum_{i=0}^{\ell-1} k_i + L_W(\mathbf{s}_\ell) \quad \text{by Lemma A.2} \\ &> (aY)^2 \sum_{i=0}^{\ell-1} z^i \\ &= (aY)^2 \frac{z^\ell - 1}{z - 1}. \end{aligned}$$

Solving for  $\ell$  finally yields the lemma.  $\square$

*Lemma A.4:* The total loss on **G1** on the last loop  $\ell$  entered is at most

$$L_W(\mathbf{s}_\ell) + (2az^{\ell/2} + 5)Y^2.$$

**Proof.** By construction of **G1**, the total loss  $L_\ell$  of **G1** on loop  $\ell$  is the total loss of  $\mathbf{GD}_{G(k_\ell, W, X)/X^2}$  on  $\mathbf{s}_\ell$ . If  $L_W(\mathbf{s}_\ell) \leq k_\ell$ , then by Theorem IV.3

$$\begin{aligned} L_\ell &\leq L_W(\mathbf{s}_\ell) + 2WX\sqrt{k_\ell} + (WX)^2 \\ &\leq L_W(\mathbf{s}_\ell) + 2Y\sqrt{k_\ell} + Y^2 \quad \text{since } Y = WX \\ &= L_W(\mathbf{s}_\ell) + (2az^{\ell/2} + 1)Y^2 \\ &< L_W(\mathbf{s}_\ell) + (2az^{\ell/2} + 5)Y^2. \end{aligned}$$

On the other hand, if  $L_W(\mathbf{s}_\ell) > k_\ell$ , then by Lemma A.1

$$\begin{aligned} L_\ell &\leq k_\ell + (2az^{\ell/2} + 5)Y^2 \\ &< L_W(\mathbf{s}_\ell) + (2az^{\ell/2} + 5)Y^2 \end{aligned}$$

and the proof is concluded.  $\square$

*Lemma A.5:* For all  $x \geq 0$ ,  $\frac{\ln(1+x)}{\ln(2.618)} \leq 0.8362\sqrt{x}$ .

**Proof.** The inequality in the statement of the lemma is equivalent to

$$\frac{\ln(1+x)}{\sqrt{x}} - 0.8362\ln(2.618) \leq 0.$$

The function  $\frac{\ln(1+x)}{\sqrt{x}}$  has a unique maximum at  $x \cong 3.921$ . At this value of  $x$  the above inequality is seen to hold.  $\square$

**Proof of Theorem IV.4.** By Lemmas A.1 and A.4,

$$\begin{aligned} &\sum_{t=1}^m (\hat{y}_t - y_t)^2 \\ &\leq \sum_{i=0}^{\ell-1} [k_i + (2az^{i/2} + 5)Y^2] \\ &\quad + L_W(\mathbf{s}_\ell) + (2az^{\ell/2} + 5)Y^2 \\ &\leq \sum_{i=0}^{\ell-1} k_i + 2aY^2 \sum_{i=0}^{\ell-1} z^{i/2} + L_W(\mathbf{s}_\ell) + 5(\ell+1)Y^2 \\ &< \sum_{i=0}^{\ell-1} L_W(\mathbf{s}_i) + L_W(\mathbf{s}_\ell) + 2aY^2 \frac{z^{(\ell+1)/2} - 1}{\sqrt{z}-1} \\ &\quad + 5(\ell+1)Y^2 \quad \text{by Lemma A.2} \\ &\leq L_W(\mathbf{s}) + 2aY^2 \frac{z^{(\ell+1)/2}}{\sqrt{z}-1} - \frac{2aY^2}{\sqrt{z}-1} + 5(\ell+1)Y^2 \\ &\leq L_W(\mathbf{s}) + 2aY^2 \frac{\sqrt{z \left( 1 + \frac{(z-1)L_W(\mathbf{s})}{(aY)^2} \right)}}{\sqrt{z}-1} - \frac{2aY^2}{\sqrt{z}-1} \\ &\quad + 5(\ell+1)Y^2 \quad \text{by Lemma A.3} \\ &= L_W(\mathbf{s}) + 2aY^2 \frac{\sqrt{z}}{\sqrt{z}-1} \sqrt{1 + \frac{(z-1)L_W(\mathbf{s})}{(aY)^2}} \\ &\quad - \frac{2aY^2}{\sqrt{z}-1} + 5(\ell+1)Y^2 \\ &< L_W(\mathbf{s}) + 2aY^2 \frac{\sqrt{z}}{\sqrt{z}-1} \left( 1 + \sqrt{\frac{(z-1)L_W(\mathbf{s})}{(aY)^2}} \right) \\ &\quad - \frac{2aY^2}{\sqrt{z}-1} + 5(\ell+1)Y^2 \\ &= L_W(\mathbf{s}) + 2aY^2 \frac{\sqrt{z}}{\sqrt{z}-1} + 2Y \frac{\sqrt{z}}{\sqrt{z}-1} \sqrt{(z-1)L_W(\mathbf{s})} \\ &\quad - \frac{2aY^2}{\sqrt{z}-1} + 5(\ell+1)Y^2 \\ &\leq L_W(\mathbf{s}) + 2Y \frac{\sqrt{z}}{\sqrt{z}-1} \sqrt{(z-1)L_W(\mathbf{s})} + 2aY^2 \\ &\quad + 5 \left[ \log_z \left( 1 + \frac{(z-1)L_W(\mathbf{s})}{(aY)^2} \right) + 1 \right] Y^2 \end{aligned}$$

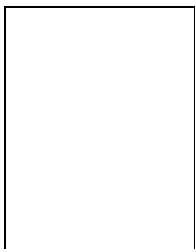
by Lemma A.3. The factor  $\frac{\sqrt{z}}{\sqrt{z}-1} \sqrt{z-1}$  is minimized at  $z \cong 2.618$ . Plugging back in this value and using Lemma A.5 we get

$$\begin{aligned} &\sum_{t=1}^m (\hat{y}_t - y_t)^2 \\ &\leq L_W(\mathbf{s}) + 6.6604Y \sqrt{L_W(\mathbf{s})} + 2aY^2 \\ &\quad + 5 \left[ 0.8362 \sqrt{\frac{1.618L_W(\mathbf{s})}{(aY)^2}} + 1 \right] Y^2 \\ &\leq L_W(\mathbf{s}) + (6.6604 + \frac{5.319}{a}) Y \sqrt{L_W(\mathbf{s})} \\ &\quad + (2a + 5)Y^2. \end{aligned}$$

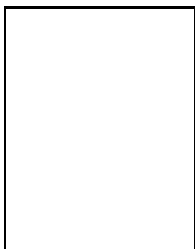
Finally, by letting  $a = 2.0979$  to trade-off between the last two terms we obtain

$$\sum_{t=1}^m (\hat{y}_t - y_t)^2 \leq L_W(\mathbf{s}) + 9.2 \left( Y \sqrt{L_W(\mathbf{s})} + Y^2 \right).$$

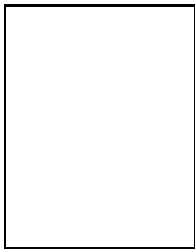
$\square$



**Nicolò Cesa Bianchi** received his M.S. (1988) and Ph.D. (1993) in information sciences from the University of Milan, Italy. He was a visiting scholar at the University of California at Santa Cruz in 1991-92 and at the Graz University of Technology, Austria, in 1994. Currently, he is an assistant professor at the department of computer science, University of Milan, Italy. His main research interests are computational learning theory and on-line analysis of algorithms.



**Phil Long** got a B.A. from Oberlin College in 1987, and a Ph.D. from the University of California at Santa Cruz in 1992. He was a postdoc at the Graz University of Technology during 1992-93, and at Duke University during 1993-95. He was at Research Triangle Institute during 1995, and is now doing another postdoc at Duke. His research interests are computational learning theory and data compression.



**Manfred Warmuth** did a Vordiplom in Informatik at Friedrich Alexander Universität. He then received a Fulbright Scholarship, and did his M.S. and Ph.D. at the University of Colorado at Boulder. He is currently a professor in the Computer Science Department of the University of California at Santa Cruz. His research interests include machine learning, on-line learning algorithms, parallel and distributed algorithms, and complexity theory.