

<https://helda.helsinki.fi>

The p -Norm Generalization of the LMS Algorithm for Adaptive Filtering

Kivinen, Jyrki

2006

Kivinen, J., Warmuth, M. K. & Hassibi, B. 2006, 'The p -Norm Generalization of the LMS Algorithm for Adaptive Filtering', IEEE Transactions on Signal Processing, vol. 54, no. 5, pp. 1782-1793. <https://doi.org/10.1109/TSP.2006.872551>

<http://hdl.handle.net/10138/27989>

<https://doi.org/10.1109/TSP.2006.872551>

publishedVersion

Downloaded from Helda, University of Helsinki institutional repository.

This is an electronic reprint of the original article.

This reprint may differ from the original in pagination and typographic detail.

Please cite the original version.

The p -Norm Generalization of the LMS Algorithm for Adaptive Filtering

Jyrki Kivinen, Manfred K. Warmuth, and Babak Hassibi

Abstract—Recently much work has been done analyzing online machine learning algorithms in a worst case setting, where no probabilistic assumptions are made about the data. This is analogous to the H^∞ setting used in adaptive linear filtering. Bregman divergences have become a standard tool for analyzing online machine learning algorithms. Using these divergences, we motivate a generalization of the least mean squared (LMS) algorithm. The loss bounds for these so-called p -norm algorithms involve other norms than the standard 2-norm. The bounds can be significantly better if a large proportion of the input variables are irrelevant, i.e., if the weight vector we are trying to learn is sparse. We also prove results for nonstationary targets. We only know how to apply kernel methods to the standard LMS algorithm (i.e., $p = 2$). However, even in the general p -norm case, we can handle generalized linear models where the output of the system is a linear function combined with a nonlinear transfer function (e.g., the logistic sigmoid).

Index Terms—Adaptive filtering, Bregman divergences, H^∞ optimality, least mean squares, online learning.

I. INTRODUCTION

WE focus on the following linear model of adaptive filtering:

$$y_t = \mathbf{u} \cdot \mathbf{x}_t + v_t. \quad (1)$$

Here \mathbf{u} is the unknown target, \mathbf{x}_t is a known input, v_t is unknown noise, and y_t is the known output signal. We are interested in algorithms that maintain a weight vector \mathbf{w}_t based on the past examples $(\mathbf{x}_\tau, y_\tau)$, $\tau = 1, \dots, t$, and, over a sequence of T trials, get as close as possible to the target \mathbf{u} . As we shall see, closely related online problems have also been studied in machine learning.

More specifically, at trial t the algorithm receives \mathbf{x}_t and y_t (in order) and has to commit to a weight vector at some point after seeing \mathbf{x}_t . We consider three problems depending on whether the algorithm needs to commit to its weight vector before or

Manuscript received December 1, 2004; revised June 26, 2005. This work was supported by the National Science Foundation under Grant CCR 9821087, the Australian Research Council, the Academy of Finland under Decision 210796, and the IST Programme of the European Community under PASCAL Network of Excellence IST-2002-506778. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Dominic K. C. Ho.

J. Kivinen is with the Department of Computer Science, University of Helsinki, FI-00014 Helsinki, Finland (e-mail: jyrki.kivinen@cs.helsinki.fi).

M. K. Warmuth is with the Department of Computer Science, University of California—Santa Cruz, Santa Cruz, CA 95064 USA (e-mail: manfred@cse.ucsc.edu).

B. Hassibi is with the Department of Electrical Engineering, California Institute of Technology, Pasadena, CA 91125 USA (e-mail: hassibi@systems.caltech.edu).

Digital Object Identifier 10.1109/TSP.2006.872551

after seeing y_t and depending on how the loss of the algorithm is measured.

- *A priori filtering*: Here we are interested in predicting the noncorrupted output $\mathbf{u} \cdot \mathbf{x}_t$ before the signal y_t is received. Therefore the algorithm needs to commit to its weight vector \mathbf{w}_{t-1} right before seeing y_t and our loss is the energy of the *a priori* filtering error $\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t$, i.e.,

$$\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2. \quad (2)$$

- *A posteriori filtering*: Here we assume that for estimating the noncorrupted output $\mathbf{u} \cdot \mathbf{x}_t$, we also have access to the measurement y_t . Thus, the algorithm needs to commit to its weight vector \mathbf{w}_t only after seeing y_t and the loss is the square of the *a posteriori* error

$$\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t)^2. \quad (3)$$

Note that as in *a priori* filtering, the algorithm does not know \mathbf{u} when it produces weight vector at trial t . It only knows the past instances and outputs.

- *Prediction*: Here we are interested in predicting the next observation y_t before receiving it. Thus the algorithm needs to commit to its weight vector \mathbf{w}_{t-1} before seeing y_t . The prediction error is $y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t$ and the loss

$$\sum_{t=1}^T (y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2. \quad (4)$$

The prediction problem of minimizing (4) is also studied in machine learning. Note that in the filtering problems, the term $v_t = y_t - \mathbf{u} \cdot \mathbf{x}_t$ is regarded as a disturbance, so we are interested in estimating the “true output” $\mathbf{u} \cdot \mathbf{x}_t$ of the linear system for the input \mathbf{x}_t . In the prediction problem we consider the y_t as the “true outcome” of some event we are interested in predicting. In that case there is no particular value in matching the prediction $\mathbf{u} \cdot \mathbf{x}_t$ at those times when it is inaccurate.

We could also define the notion of *a posteriori* prediction, i.e., trying to minimize

$$\sum_{t=1}^T (y_t - \mathbf{w}_t \cdot \mathbf{x}_t)^2. \quad (5)$$

However, since y_t is known when \mathbf{w}_t is chosen, the loss (5) is trivially minimized by just choosing \mathbf{w}_t such that $\mathbf{w}_t \cdot \mathbf{x}_t = y_t$.

Although there are algorithms that do satisfy $\mathbf{w}_t \cdot \mathbf{x}_t = y_t$ in some limiting cases, taking this condition as the primary design principle does not seem to add anything. Hence, we do not further consider the loss (5).

In contrast to the loss function used by the prediction problem, the loss functions for the two filtering problems include the target \mathbf{u} that is unknown. Because the algorithm cannot even evaluate its own loss, we need to be careful about setting a reasonable performance criterion. We next set the performance criteria we use in this paper, starting with *a priori* filtering and its connection to recent work in machine learning.

Clearly the quality of output depends on the amount of noise, which can be defined, for example, as $\sum_{t=1}^T (y_t - \mathbf{u} \cdot \mathbf{x}_t)^2$. Additionally, even with no noise, the loss (2) for any given algorithm can be made arbitrarily large by scaling \mathbf{u} . To have a well-defined choice of \mathbf{u} , we consider the *regularized loss* $\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + (1/\eta)\|\mathbf{u}\|_2^2$ where $\eta > 0$ is a tradeoff parameter. We then normalize the algorithm's loss (2) with respect to the regularized loss. Since we wish to avoid assumptions about \mathbf{u} , we consider the worst case choice, leading us to the quantity

$$\max_{\mathbf{u}} \frac{\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2}{\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + \frac{1}{\eta}\|\mathbf{u}\|_2^2}. \quad (6)$$

Given the data (\mathbf{x}_t, y_t) and an algorithm for producing \mathbf{w}_t , the quantity (6) is always well defined. In control theory, (6) is seen as a maximum energy gain and called the H^∞ norm. (For the above, and as done throughout this paper, we assumed $\mathbf{w}_0 = \mathbf{0}$; if $\mathbf{w}_0 \neq \mathbf{0}$, then $\|\mathbf{u}\|_2^2$ must be replaced by $\|\mathbf{u} - \mathbf{w}_0\|_2^2$.)

To get a reference point, consider the least mean squares (LMS) algorithm [2] (also known as the Widrow–Hoff algorithm), defined by the update rule

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - y_t)\mathbf{x}_t \quad (7)$$

where $\eta > 0$ is now a parameter of the algorithm and called the learning rate. According to the basic result for *a priori* filtering [3], if $\eta \leq 1/\max\|\mathbf{x}_t\|_2^2$, then the LMS algorithm satisfies

$$\max_{\mathbf{u}} \frac{\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2}{\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + \frac{1}{\eta}\|\mathbf{u}\|_2^2} \leq 1. \quad (8)$$

In other words, LMS has H^∞ norm at most 1. (Notice that the learning rate parameter of the algorithm becomes the tradeoff parameter for the regularized loss.) Further, no algorithm can have H^∞ norm less than 1. Therefore, we say that LMS is H^∞ optimal.

To compare this with results from machine learning, assume there is a known upper bound X_2 such that $\|\mathbf{x}_t\|_2 \leq X_2$ for all t , and write $\eta = \alpha/X_2^2$. Then Cesa–Bianchi *et al.* [4] have shown that for $0 < \alpha < 1$

$$\sum_{t=1}^T (y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 \leq \frac{1}{1-\alpha} \sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + \frac{1}{\alpha} X_2^2 \|\mathbf{u}\|_2^2. \quad (9)$$

To compare prediction with filtering, we write (6) as

$$\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 \leq \sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + \frac{1}{\alpha} X_2^2 \|\mathbf{u}\|_2^2 \quad (10)$$

where X_2 and η are as above and $0 < \alpha \leq 1$. We see that the bounds are similar in form, except for the factor $1/(1-\alpha)$ in (9).

The factor $1/(1-\alpha)$ in (9) is a source of many difficulties in machine learning, where the goal is to tune the learning rate so as to obtain the smallest possible bound. However, the filtering bound (10) is optimized at $\alpha = 1$. Thus we omit the α parameter from the filtering bounds when the norm of instances is bounded.

Motivated by the similarity between (9) and (10), we are going to take machine learning techniques that have recently been used to generalize the LMS algorithm and apply them in the filtering setting. This leads to generalizations of (10) and new interpretations of the filtering algorithms. Techniques we are interested in include:

- 1) motivating algorithms in terms of minimization problems based on Bregman divergences [5], [6];
- 2) replacing the 2-norms in the bounds by other norms [5], [7], [8];
- 3) allowing for nonstationary targets [9] and nonlinear predictors [10].

Before going on with the above program, let us have a brief look at the *a posteriori* model. The H^∞ norm for *a posteriori* filtering is

$$\max_{\mathbf{u}} \frac{\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t)^2}{\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + \frac{1}{\eta}\|\mathbf{u}\|_2^2}.$$

Notice that since y_t is available when choosing \mathbf{w}_t , we can trivially obtain H^∞ norm at most 1 by any choice that satisfies $\mathbf{w}_t \cdot \mathbf{x}_t = y_t$. One particular way of doing this would be to let the learning rate go to infinity in the *normalized LMS* algorithm [3]. However, there are other criteria that are minimized by using a finite learning rate, while still retaining the H^∞ norm at most 1. For example, this is the case if the data points are generated by the model (1) with the noise variables v_t independent and Gaussian [3, Theorem 9]. Thus, while requiring the H^∞ norm to be at most 1 is a good robustness guarantee, in the *a posteriori* case such a worst case measure is not by itself a sufficient criterion for choosing a good algorithm. In the following we will state all our bounds both for *a priori* and *a posteriori* filtering, but they must be read with this caveat in mind.

Our H^∞ -based performance criteria do not directly address convergence. If the data are generated by the model (1) with the noise variables v_i independent and Gaussian, then one could hope that the weights \mathbf{w}_t would converge toward the target \mathbf{u} . However, if we do not wish to make such assumptions about noise, the issue becomes less clear. An algorithm geared toward fast convergence under zero-mean independent noise may fail badly if, say, the early data points have large amounts of biased and correlated noise. We aim for results that are not sensitive to probabilistic assumptions and develop bounds like (6) and (10), which hold for every sequence of examples. Such worst

case bounds are rather stringent. If the examples are independent identically distributed (i.i.d.), an averaging technique can be used to convert worst case loss bounds to bounds on the expected loss (see, e.g., [5, Section 8]) or bounds on the probability of high loss [11]. Clearly the choice of algorithm should depend on the assumptions. In particular, even with independent noise, updates like (7) with fixed learning rate do not typically lead to convergence but remain oscillating around the optimal weight setting.

In Section III, we introduce Bregman divergences and show how a Bregman divergence can be used to derive two subtly different updates: the *implicit* and *explicit* update. When the squared Euclidean distance is used as the Bregman divergence, these updates give the standard *LMS* and *normalized LMS* algorithm [3], respectively. In Section IV, we give filtering loss bounds for the explicit and implicit updates in the case of Bregman divergences based on squared q -norms [7]. These bounds generalize the results of Hassibi *et al.* [3] about the H^∞ optimality of LMS and normalized LMS for the *a priori* and *a posteriori* filtering problems. The generalization replaces the product $\|\mathbf{x}\|_2\|\mathbf{u}\|_2$ in the bound by another product of dual norms $\|\mathbf{x}\|_p\|\mathbf{u}\|_q$, where p and q are such that $1/p + 1/q = 1$ and $2 \leq p < \infty$. The new bounds are significantly stronger when the target \mathbf{u} is sparse, i.e., has few nonzero components. In Section V, we generalize the q -norm based algorithms to allow for nonstationary targets \mathbf{u}_t . The loss bounds in the nonstationary case include an extra term that depends on the total distance \mathbf{u}_t travels during the whole sequence, as measured by the q -norm. Again there are no distribution assumptions about this movement. Section VI gives bounds for generalized linear regression where the linear predictor is fed through a nonlinear transfer function (such as the logistic sigmoid). Some simulations are reported in Section VII, and our conclusions presented in Section VIII.

Some preliminary results of this paper were presented at the 13th IFAC Symposium on System Identification [1]. This paper includes some additional algorithms and new simulation results, as well as full proofs of the theoretical results.

II. THE LMS BOUND

As an introduction to our methods, we rederive the basic result of [3]. Later we will see how the algorithm and proof generalize from the Euclidean to other p -norms.

Theorem 1 [3]: Assume that $\|\mathbf{x}_t\|_2 \leq X_2$ for all t , and choose $\eta = 1/X_2^2$. Then the LMS algorithm (7) satisfies

$$\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 \leq \sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + X_2^2 \|\mathbf{u}\|_2^2$$

for any $\mathbf{u} \in \mathbf{R}^n$.

Proof: Following [4], we analyze the *progress* $d_t = (1/2)\|\mathbf{u} - \mathbf{w}_t\|_2^2 - (1/2)\|\mathbf{u} - \mathbf{w}_{t-1}\|_2^2$ made at update t toward the *comparison vector* \mathbf{u} . Direct calculation gives us

$$d_t = \eta(y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)(\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t) - \frac{\eta^2}{2}(y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 \|\mathbf{x}_t\|_2^2.$$

By estimating $\|\mathbf{x}_t\|_2 \leq X_2$ and rearranging terms, we get

$$d_t \geq \frac{\eta}{2}s_t^2 - \frac{\eta}{2}r_t^2 + \frac{\eta}{2}(s_t - r_t)^2(1 - \eta X_2^2)$$

where $s_t = \mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t$ and $r_t = \mathbf{u} \cdot \mathbf{x}_t - y_t$. Since $\eta X_2^2 = 1$ and $\mathbf{w}_0 = \mathbf{0}$, we can apply $\|\mathbf{u} - \mathbf{w}_0\|_2 = \|\mathbf{u}\|_2$ and $\|\mathbf{u} - \mathbf{w}_{T+1}\|_2 \geq 0$ to get

$$\begin{aligned} \frac{1}{2}\|\mathbf{u}\|_2^2 &\geq \frac{1}{2}\|\mathbf{u} - \mathbf{w}_0\|_2^2 - \frac{1}{2}\|\mathbf{u} - \mathbf{w}_{T+1}\|_2^2 \\ &= \sum_{t=1}^T d_t \\ &\geq \frac{1}{2X_2^2} \left(\sum_{t=1}^T s_t^2 - \sum_{t=1}^T r_t^2 \right) \end{aligned}$$

from which the claim follows. \blacksquare

III. DERIVATION OF ALGORITHMS

In this section we give the basic definitions of Bregman divergences and explain their use in deriving generalizations of the LMS algorithm. (See [12] and references therein for more background on these divergences.) Later the same Bregman divergences will be used to prove bounds for these new algorithms. Note that the bound for the LMS algorithm involves the 2-norms of the inputs \mathbf{x} and target \mathbf{u} . The bounds for the new algorithm will depend on norms $\|\mathbf{x}\|_p$ and $\|\mathbf{u}\|_q$ where in general $p, q \neq 2$.

Assume that F is a strictly convex twice differentiable function from a subset of \mathbf{R}^n to \mathbf{R} . Denote its gradient by $\mathbf{f} = \nabla F$; notice that \mathbf{f} is one-to-one. The *Bregman divergence* $\Delta_F(\mathbf{u}, \mathbf{w})$ [13] is defined for $\mathbf{u}, \mathbf{w} \in \mathbf{R}^n$ as the error in approximating $F(\mathbf{u})$ by its first order Taylor polynomial around \mathbf{w} . More formally

$$\Delta_F(\mathbf{u}, \mathbf{w}) = F(\mathbf{u}) - F(\mathbf{w}) - (\mathbf{u} - \mathbf{w}) \cdot \mathbf{f}(\mathbf{w}).$$

The Bregman divergence $\Delta_F(\mathbf{u}, \mathbf{w})$ is always nonnegative, and zero only for $\mathbf{u} = \mathbf{w}$. It is (strictly) convex in \mathbf{u} but might not be convex in \mathbf{w} . Usually, Δ_F is not symmetric.

Example 1: For $q > 1$, define $F(\mathbf{w}) = (1/2)\|\mathbf{w}\|_q^2$, where $\|\cdot\|_q$ denotes the q -norm defined as $\|\mathbf{w}\|_q = (\sum_i |w_i|^q)^{1/q}$. We denote the corresponding Bregman divergence by Δ_q . Thus

$$\Delta_q(\mathbf{u}, \mathbf{w}) = \frac{1}{2}\|\mathbf{u}\|_q^2 - \frac{1}{2}\|\mathbf{w}\|_q^2 - (\mathbf{u} - \mathbf{w}) \cdot \mathbf{f}(\mathbf{w})$$

where the gradient is given by

$$f_i(\mathbf{w}) = \frac{\text{sign}(w_i)|w_i|^{q-1}}{\|\mathbf{w}\|_q^{q-2}}.$$

A second important family of Bregman divergences is the relative entropy and its variants.

Example 2: Assume $w_i \geq 0$ for all i and define $F(\mathbf{w}) = \sum_i (w_i \ln w_i - w_i)$, with the usual convention $0 \ln 0 = 0$. Then

$$\Delta_F(\mathbf{u}, \mathbf{w}) = \sum_i \left(u_i \ln \frac{u_i}{w_i} - u_i + w_i \right)$$

is the unnormalized relative entropy. (When $\sum_i u_i = \sum_i w_i = 1$, this gives the standard relative entropy.) The gradient is given by $f_i(\mathbf{w}) = \ln w_i$.

The following generalization of the Pythagorean theorem follows directly from the definition of a Bregman divergence:

$$\Delta_F(\mathbf{u}, \mathbf{w}') = \Delta_F(\mathbf{u}, \mathbf{w}) + \Delta_F(\mathbf{w}, \mathbf{w}') + (\mathbf{w} - \mathbf{u}) \cdot (\mathbf{f}(\mathbf{w}') - \mathbf{f}(\mathbf{w})). \quad (11)$$

Since the dot product $(\mathbf{w} - \mathbf{u}) \cdot (\mathbf{f}(\mathbf{w}') - \mathbf{f}(\mathbf{w}))$ can be positive, this shows in particular that Δ_F does not satisfy the triangle inequality. We recover the standard Pythagorean theorem when the divergence is the squared Euclidean distance (i.e., \mathbf{f} is identity) and the dot product is zero (i.e., $(\mathbf{w}' - \mathbf{w})$ and $\mathbf{w} - \mathbf{u}$ are orthogonal).

We now use a Bregman divergence Δ_F as a regularizer for deriving an update rule. This framework for motivating updates was introduced in [5] in the prediction setting. In the following, we are mainly interested in Bregman divergences based on the squared q -norm. They were introduced in [7] to analyze algorithms for learning linear threshold functions.

Suppose an example (\mathbf{x}_t, y_t) has been observed and we wish to update our hypothesis \mathbf{w}_{t-1} based on this example. We wish to decrease the squared loss $(y_t - \mathbf{w} \cdot \mathbf{x}_t)^2$ (other convex loss functions can also be considered; see Section VI). However, we should not make big changes based on just a single example. Thus, we define

$$C_t(\mathbf{w}) = \Delta_F(\mathbf{w}, \mathbf{w}_{t-1}) + \frac{1}{2}\eta(y_t - \mathbf{w} \cdot \mathbf{x}_t)^2$$

where $\eta > 0$ is a tradeoff parameter, and tentatively set $\mathbf{w}_t = \arg \min_{\mathbf{w}} C_t(\mathbf{w})$. Since C_t is convex, we can minimize by setting $\nabla C_t(\mathbf{w}_t) = 0$. By substituting the definition of Δ_F , this becomes

$$\mathbf{w}_t = \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}_{t-1}) - \eta(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - y_t)\mathbf{x}_t). \quad (12)$$

Since \mathbf{w}_t appears on both sides of (12), we call the update rule defined by this equality the *implicit update* for divergence Δ_F . Notice that (12) can be solved numerically by a line search since $\mathbf{w} = \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}_{t-1} + \alpha\mathbf{x}_t))$ for some scalar α , and the inverse \mathbf{f}^{-1} is easy to compute in the cases we consider. Also in the special case of 2-norm ($\Delta_F = \Delta_2$), with \mathbf{f} the identity function, we can solve (12) in closed form to get

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\eta}{1 + \eta\|\mathbf{x}_t\|_2^2}(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - y_t)\mathbf{x}_t. \quad (13)$$

This is the algorithm called *normalized LMS* in [3].

Instead of solving (12) numerically, we often find it sufficient to notice that for reasonable values of η , the values $\mathbf{w}_t \cdot \mathbf{x}_t$ and $\mathbf{w}_{t-1} \cdot \mathbf{x}_t$ should be fairly close to each other. Thus, we may approximate the solution of (12) by

$$\mathbf{w}_t = \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}_{t-1}) - \eta(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - y_t)\mathbf{x}_t). \quad (14)$$

We call this the *explicit update* for divergence Δ_F . The special case $\Delta_F = \Delta_2$ gives the usual LMS algorithm.

Note that the explicit update uses the gradient of the square loss evaluated at the old weight vector \mathbf{w}_{t-1} , whereas the implicit update is based on the gradient at the updated parameter

vector \mathbf{w}_t . For a discussion of taking the old gradient versus the future gradient in for the prediction problem, and a derivation of the implicit LMS algorithm, see [5]. In [14], an implicit update was derived as an alternate to the TD(λ) algorithm. In this case the implicit definition was crucial for producing an improved algorithm.

IV. BOUNDS IN TERMS OF DIFFERENT NORMS

Our interest in considering the generalization of LMS to the p -norm based algorithms comes from the fact that for these algorithms, the term $\|\mathbf{x}\|_2\|\mathbf{u}\|_2$ in the LMS bound is replaced by another product of dual norms $\|\mathbf{x}\|_p\|\mathbf{u}\|_q$ (i.e., $1/p + 1/q = 1$). We discuss the implications of this after giving the main result, which is a direct generalization of Theorem 1.

We consider the explicit (14) and implicit (12) updates for the divergence $\Delta_q(\mathbf{u}, \mathbf{w})$ given in Example 1. The special case $q = 2$ gives the classic LMS and Theorem 1. For the updates, we need the gradient \mathbf{f} , which was given in Example 1, and also its inverse \mathbf{f}^{-1} , which is easily seen to be

$$f_i^{-1}(\theta) = \frac{\text{sign}(\theta_i)|\theta_i|^{p-1}}{\|\boldsymbol{\theta}\|_2^{p-2}}$$

where $1/p + 1/q = 1$.

We assume the relationship $1/p + 1/q = 1$ throughout this paper. It means that we can apply *Hölder's inequality* $|\mathbf{w} \cdot \mathbf{x}| \leq \|\mathbf{w}\|_q\|\mathbf{x}\|_p$. As a further convention, we assume $q \leq p$, so $1 < q \leq 2 \leq p < \infty$. The important special case $p = q = 2$ gives $\Delta_2(\mathbf{u}, \mathbf{w}) = (1/2)\|\mathbf{u} - \mathbf{w}\|_2^2$, with \mathbf{f} the identity function.

We use the following inequality for proving bounds for the updates:

$$\Delta_q(\mathbf{w}, \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}) + \mathbf{x})) \leq \frac{p-1}{2}\|\mathbf{x}\|_p^2. \quad (15)$$

This inequality is implied by derivations given in [7] and was stated explicitly in [8, Lemma 2]. For completeness, we give the proof in Appendix I.

Theorem 2: Fix p and q such that $1/p + 1/q = 1$ and $2 \leq p < \infty$. Assume that $\|\mathbf{x}_t\|_p \leq X_p$ for all t . Then the explicit update (14) for Δ_q with learning rate $\eta = 1/((p-1)X_p^2)$ satisfies

$$\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 \leq \sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + (p-1)X_p^2\|\mathbf{u}\|_q^2$$

for any $\mathbf{u} \in \mathbb{R}^n$.

Proof: Following [5], we analyze the *progress* $d_t = \Delta_q(\mathbf{u}, \mathbf{w}_{t-1}) - \Delta_q(\mathbf{u}, \mathbf{w}_t)$ made at update t toward the *comparison vector* \mathbf{u} . By substituting (14) into (11) and then using (15), we get

$$\begin{aligned} d_t &= \eta(y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)\mathbf{x}_t \cdot (\mathbf{u}_t - \mathbf{w}_{t-1}) - \Delta_q(\mathbf{w}_{t-1}, \mathbf{w}_t) \\ &\geq \eta(y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)(\mathbf{u}_t \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t) \\ &\quad - \frac{p-1}{2}\eta^2(y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 X_p^2. \end{aligned}$$

By rearranging terms, we can write this as

$$d_t \geq \frac{\eta}{2}s_t^2 - \frac{\eta}{2}r_t^2 + \frac{\eta}{2}(s_t - r_t)^2 (1 - \eta(p-1)X_p^2)$$

where $s_t = \mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t$ and $r_t = \mathbf{u} \cdot \mathbf{x}_t - y_t$. Since $\eta(p-1)X_p^2 = 1$ and $\mathbf{w}_0 = \mathbf{0}$, we can apply $\Delta_q(\mathbf{u}, \mathbf{w}_0) = (1/2)\|\mathbf{u}\|_q^2$ and $\Delta_q(\mathbf{u}, \mathbf{w}_{T+1}) \geq 0$ to get

$$\begin{aligned} \frac{\|\mathbf{u}\|_q^2}{2} &\geq \Delta_q(\mathbf{u}, \mathbf{w}_0) - \Delta_q(\mathbf{u}, \mathbf{w}_{T+1}) \\ &= \sum_{t=1}^T d_t \\ &\geq \frac{1}{2(p-1)X_p^2} \left(\sum_{t=1}^T s_t^2 - \sum_{t=1}^T r_t^2 \right) \end{aligned}$$

from which the claim follows. \blacksquare

The main intuitive implication of Theorem 2 (and later Theorem 3, which will deal with the implicit update) is that the bound favors large p when the target \mathbf{u} is sparse. To make this more precise, we compare the bound for $p = 2$ (i.e., classic LMS) against $p = 2 \ln n$ (i.e., fairly large p). Gentile and Littlestone [8, Corollary 7] have shown that for the particular choice $p = 2 \ln n$, we have

$$(p-1)\|\mathbf{x}\|_p^2 \|\mathbf{u}\|_q^2 \leq (2e \ln n)\|\mathbf{x}\|_\infty^2 \|\mathbf{u}\|_1^2 \quad (16)$$

(where $\|\mathbf{x}\|_\infty = \max_i |x_i|$). Thus, we compare the bound $\|\mathbf{x}\|_2^2 \|\mathbf{u}\|_2^2$ (for LMS) with the bound $(2e \ln n)\|\mathbf{x}\|_\infty^2 \|\mathbf{u}\|_1^2$ (for large p).

Since the p -norm is decreasing in p , we have $\|\mathbf{u}\|_2 \leq \|\mathbf{u}\|_1$ and $\|\mathbf{x}\|_2 \geq \|\mathbf{x}\|_\infty$, with equality if the vector has only one nonzero component. Hence, the dependence on \mathbf{u} favors $p = 2$, but the advantage gets smaller if \mathbf{u} is very sparse. Similarly, the dependence on \mathbf{x} favors large p , but the advantage gets smaller if \mathbf{x} is very sparse.

To get a concrete picture of the tradeoff, let us consider two extreme cases. In the first case, we choose $\mathbf{u} \in \{-1, 1\}^n$ and $\mathbf{x} \in \{-1, 0, 1\}^n$ such that exactly one component x_i is nonzero. Then $\|\mathbf{u}\|_2^2 = n$, $\|\mathbf{u}\|_1^2 = n^2$, and $\|\mathbf{x}\|_2 = \|\mathbf{x}\|_\infty = 1$. The LMS bound becomes simply n , while the large p bound becomes $2en^2 \ln n$. Hence, the LMS bound is clearly better for large n . In the second case, choose $\mathbf{u} \in \{-1, 0, 1\}^n$ such that exactly one component u_i is nonzero, and choose $\mathbf{x} \in \{-1, 1\}^n$. Then $\|\mathbf{u}\|_2 = \|\mathbf{u}\|_1 = 1$, $\|\mathbf{x}\|_2^2 = n$, and $\|\mathbf{x}\|_\infty = 1$. The LMS bound is n as in the first case, but the large p bound drops to $2e \ln n$. Notice that the dependence on n in this last bound is only logarithmic, so for large p the difference to LMS can be quite large.

The above two example scenarios were of course unrealistically extreme. In a typical application, one would expect the components x_i of the inputs \mathbf{x} to have roughly the same magnitude, so the inputs would be relatively dense. Then a large p would be favored if $\|\mathbf{u}\|_1$ is close to $\|\mathbf{u}\|_2$, which is the case if most of the weight in \mathbf{u} is concentrated on only few components. One should also notice that the upper bounds might not reflect the actual behavior of the algorithms. However, simulations suggest that the picture given here is at least qualitatively correct: the algorithms for $p = 2$ and large p are incomparable, and large p is better if the target is sparse. See Section VII for some examples.

In the context of prediction, much attention has been paid to *multiplicative algorithms* such as Winnow [15] and EG [5], which have bounds similar to the p -norm algorithms for $p = O(\log n)$. In addition to upper bounds and simulations [5], there are also some lower bounds [16] showing that in certain situations LMS-style algorithms cannot perform as well as multiplicative ones. The multiplicative EG algorithm can be seen as applying the update (14) with $f_i(\mathbf{w}) = \ln w_i$ (with a further normalization step). The analysis of EG can also be lifted to the filtering setting, resulting in the bound

$$\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 \leq \sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + \ln(2n)\|\mathbf{x}\|_\infty^2 \|\mathbf{u}\|_1^2$$

for a scaled explicit version. See Appendix II for details and notice the improved constant of $\ln 2n$ over $2e \ln n$ appearing in (16). Multiplicative algorithms are closely related to L_1 regularization, which can be seen as a form of feature selection [17].

We now consider the *a posteriori* case. The following theorem generalizes the result about normalized LMS in [3]. However, our result has an additional restriction on the learning rate, which we believe to be an artefact of the proof technique. We shall discuss this after giving the theorem and its proof.

Theorem 3: Fix p and q such that $1/p + 1/q = 1$ and $2 \leq p < \infty$. Assume that $\|\mathbf{x}_t\|_p \leq X_p$ for all t . Then the implicit update for Δ_q with learning rate $\eta = 1/((p-1)X_p^2)$ satisfies

$$\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t)^2 \leq \sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + (p-1)X_p^2 \|\mathbf{u}\|_q^2.$$

Proof: Again let $d_t = \Delta_q(\mathbf{u}, \mathbf{w}_{t-1}) - \Delta_q(\mathbf{u}, \mathbf{w}_t)$. By substituting (12) into (11) and applying (15), we get

$$\begin{aligned} d_t &= \eta(y_t - \mathbf{w}_t \cdot \mathbf{x}_t) \mathbf{x}_t \cdot (\mathbf{u} - \mathbf{w}_{t-1}) - \Delta_q(\mathbf{w}_{t-1}, \mathbf{w}_t) \\ &\geq \eta(y_t - \mathbf{w}_t \cdot \mathbf{x}_t)(y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t) - \eta(y_t - \mathbf{w}_t \cdot \mathbf{x}_t) \\ &\quad \times (y_t - \mathbf{u} \cdot \mathbf{x}_t) - \frac{p-1}{2} \eta^2 (y_t - \mathbf{w}_t \cdot \mathbf{x}_t)^2 X_p^2. \end{aligned}$$

Since \mathbf{w}_t minimizes C_t , it is easy to show that $\mathbf{w}_{t-1} \cdot \mathbf{x}_t \leq \mathbf{w}_t \cdot \mathbf{x}_t \leq y_t$ or $y_t \leq \mathbf{w}_t \cdot \mathbf{x}_t \leq \mathbf{w}_{t-1} \cdot \mathbf{x}_t$; that is, the update moves $\mathbf{w} \cdot \mathbf{x}_t$ to the right direction but not too far. This implies

$$(y_t - \mathbf{w}_t \cdot \mathbf{x}_t)(y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t) \geq (y_t - \mathbf{w}_t \cdot \mathbf{x}_t)^2 \quad (17)$$

so we get

$$\begin{aligned} d_t &\geq \eta(y_t - \mathbf{w}_t \cdot \mathbf{x}_t)^2 - \eta(y_t - \mathbf{w}_t \cdot \mathbf{x}_t)(y_t - \mathbf{u} \cdot \mathbf{x}_t) \\ &\quad - \frac{p-1}{2} \eta^2 (y_t - \mathbf{w}_t \cdot \mathbf{x}_t)^2 X_p^2. \end{aligned}$$

We can rewrite this as

$$d_t \geq \eta(s_t - r_t)^2 + \eta(s_t - r_t)r_t - \frac{p-1}{2} \eta^2 X_p^2 (s_t - r_t)^2$$

where $s_t = \mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t$ and $r_t = \mathbf{u} \cdot \mathbf{x}_t - y_t$. By rearranging terms, this becomes

$$d_t \geq \frac{\eta}{2} s_t^2 - \frac{\eta}{2} r_t^2 + \frac{\eta}{2} (s_t - r_t)^2 (1 - \eta(p-1)X_p^2).$$

The rest follows as in the proof of Theorem 2. \blacksquare

Our proof actually implies

$$\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t)^2 \leq \sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + \frac{1}{\eta} \|\mathbf{u}\|_q^2 \quad (18)$$

for any learning rate $0 < \eta \leq 1/((p-1)X_p^2)$. For the case $p = 2$, Hassibi *et al.* [3] actually show (18) for any $\eta > 0$. Notice that the estimate (17) in our proof can equivalently be written as $(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t)/(\mathbf{w}_t \cdot \mathbf{x}_t - y_t) \geq 0$. This holds as equality for $\eta = 0$, but becomes very loose as η approaches infinity (so $\mathbf{w}_t \cdot \mathbf{x}_t - y_t$ approaches zero). In the case $p = 2$, we can use the closed form (13) of the normalized LMS algorithm to obtain $(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t)/(\mathbf{w}_t \cdot \mathbf{x}_t - y_t) = \eta \|\mathbf{x}_t\|_2^2$. Using this tighter estimate allows the proof to go through for arbitrary $\eta > 0$. Unfortunately, we have not been able to obtain a similar bound for the case $p > 2$, with nonlinear \mathbf{f} in the update (12).

As discussed in [5], whenever a learning rate η needs to be tuned, then the tuned choice should be of the correct ‘‘type.’’ As we shall see, this is indeed the case in the above two theorems. We denote the type of the weight vectors as $[\mathbf{w}]$ and the type of the instances as $[\mathbf{x}]$. The type of the outputs must then be $[\mathbf{w} \cdot \mathbf{x}] = [\mathbf{w}][\mathbf{x}]$. It is easy to check that the transformations \mathbf{f} and \mathbf{f}^{-1} for Δ_p do not change the type of a weight vector. So now the type of η in the implicit and explicit update for Δ_q must be $[\mathbf{x}]^{-2}$ and the tunings prescribed in the theorems indeed choose an η of this type. Throughout this paper, our tunings of η always fix the type of η for all the updates discussed.

V. NONSTATIONARY TARGETS

Following [9], we now consider a variant of the algorithm that keeps the q -norm of the weight vector bounded by U_q , where $U_q > 0$ is a parameter to the algorithm. We call this two-step update the *bounded explicit update* for Δ_F .

- *Explicit update step:* Let $\mathbf{w}'_t = \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}_{t-1}) - \eta(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - y_t)\mathbf{x}_t)$.
- *Out-of-bound update step:* If $\|\mathbf{w}'_t\|_q > U_q$, then $\mathbf{w}_t = U_q \mathbf{w}'_t / \|\mathbf{w}'_t\|_q$; otherwise $\mathbf{w}_t = \mathbf{w}'_t$.

Thus if the update tries to increase the q -norm of its weight vector above U_q , then we scale it back.

We now let the target \mathbf{u}_t vary with time (nonstationary model):

$$y_t = \mathbf{u}_t \cdot \mathbf{x}_t + v_t. \quad (19)$$

As previously, our bound will include a penalty for the (maximum) norm of \mathbf{u}_t . Additionally, there is now also a penalty for the total distance the target moves during the process.

Theorem 4: Fix p and q such that $1/p + 1/q = 1$ and $2 \leq p < \infty$. Assume $\|\mathbf{x}_t\|_p \leq X_p$ and $\|\mathbf{u}_t\|_q \leq U_q$ for all t . Then the bounded explicit update for Δ_q with learning rate $\eta = 1/((p-1)X_p^2)$ and parameter U_q satisfies

$$\begin{aligned} \sum_{t=1}^T (\mathbf{u}_t \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 &\leq \sum_{t=1}^T (\mathbf{u}_t \cdot \mathbf{x}_t - y_t)^2 + (p-1)X_p^2 U_q^2 \\ &\quad + 2(p-1)X_p^2 U_q \sum_{t=1}^{T-1} \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q. \end{aligned}$$

Proof: We apply the proof technique introduced in the prediction setting in [9]. We define the progress at trial t as the sum of three parts $d_t = d_t^1 + d_t^2 + d_t^3$, where

$$\begin{aligned} d_t^1 &= \Delta_q(\mathbf{u}_t, \mathbf{w}_{t-1}) - \Delta_q(\mathbf{u}_t, \mathbf{w}'_t) \\ d_t^2 &= \Delta_q(\mathbf{u}_t, \mathbf{w}'_t) - \Delta_q(\mathbf{u}_t, \mathbf{w}_t) \\ d_t^3 &= \Delta_q(\mathbf{u}_t, \mathbf{w}_t) - \Delta_q(\mathbf{u}_{t+1}, \mathbf{w}_t). \end{aligned}$$

Then $d_t = \Delta_q(\mathbf{u}_t, \mathbf{w}_{t-1}) - \Delta_q(\mathbf{u}_{t+1}, \mathbf{w}_t)$. (For notational convenience we define $\mathbf{u}_{T+1} = \mathbf{u}_T$ for the last time step.)

For $\eta \leq 1/((p-1)X_p^2)$, the proof of Theorem 2 gives directly

$$d_t^1 \geq \frac{\eta}{2} s_t^2 - \frac{\eta}{2} r_t^2 \quad (20)$$

where $s_t = \mathbf{u}_t \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t$ and $r_t = \mathbf{u}_t \cdot \mathbf{x}_t - y_t$.

For estimating d_t^2 , first note that the out-of-bound step can be expressed as

$$\mathbf{w}_t = \arg \min_{\|\mathbf{w}\|_q \leq U_q} \Delta_q(\mathbf{w}, \mathbf{w}'_t).$$

In other words, \mathbf{w}_t is the *projection* of \mathbf{w}'_t into the closed convex set $B = \{\mathbf{w} \mid \|\mathbf{w}\|_q \leq U_q\}$ with respect to Δ_q . Well-known properties of such projections [9], [13] imply that for any $\mathbf{u} \in B$, we have $\Delta_q(\mathbf{u}, \mathbf{w}_t) \leq \Delta_q(\mathbf{u}, \mathbf{w}'_t)$ and thus $d_t^2 \geq 0$.

From the definition of Δ_q , we get

$$d_t^3 = \frac{1}{2} \|\mathbf{u}_t\|_q^2 - \frac{1}{2} \|\mathbf{u}_{t+1}\|_q^2 + (\mathbf{u}_{t+1} - \mathbf{u}_t) \cdot \mathbf{f}(\mathbf{w}_t).$$

By Hölder’s inequality, $|(\mathbf{u}_{t+1} - \mathbf{u}_t) \cdot \mathbf{f}(\mathbf{w}_t)| \leq \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q \|\mathbf{f}(\mathbf{w}_t)\|_p$. Since $\|\mathbf{f}(\mathbf{w}_t)\|_p = \|\mathbf{w}_t\|_q \leq U_q$, we get

$$d_t^3 \geq \frac{1}{2} \|\mathbf{u}_t\|_q^2 - \frac{1}{2} \|\mathbf{u}_{t+1}\|_q^2 - U_q \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q.$$

By summing over $t = 1, \dots, T$ and substituting the value of η , we obtain

$$\begin{aligned} &\Delta_q(\mathbf{u}_1, \mathbf{w}_0) - \Delta_q(\mathbf{u}_{T+1}, \mathbf{w}_T) \\ &= \sum_{t=1}^T d_t \\ &\geq \frac{1}{2(p-1)X_p^2} \sum_{t=1}^T s_t^2 - \frac{1}{2(p-1)X_p^2} \sum_{t=1}^T r_t^2 \\ &\quad + \frac{1}{2} \|\mathbf{u}_1\|_q^2 - \frac{1}{2} \|\mathbf{u}_{T+1}\|_q^2 + U_q \sum_{t=1}^T \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q. \end{aligned}$$

For $\mathbf{w}_0 = \mathbf{0}$, we have $\Delta_q(\mathbf{u}_1, \mathbf{w}_0) = \|\mathbf{u}_1\|_q^2/2$. Estimating $\Delta_q(\mathbf{u}_{T+1}, \mathbf{w}_T) \geq 0$ and $\|\mathbf{u}_{T+1}\|_q \leq U_q$ gives the claim. ■

In the special case $\mathbf{u}_t = \mathbf{u}_{t+1}$ for all t , the result becomes Theorem 2 with the exception that the norm bound U_q must be fixed in advance.

The same technique can be applied to the *a posteriori* problem. Given $U_q > 0$, we define the *bounded implicit update* for Δ_F with the following two-step update.

- *Implicit update step:* Let \mathbf{w}_t be such that

$$\mathbf{w}_t = \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}'_{t-1}) - \eta(\mathbf{w}_t \cdot \mathbf{x}_t - y_t)\mathbf{x}_t).$$

- *Out-of-bound update step:* If $\|\mathbf{w}_t\|_q > U_q$, then $\mathbf{w}'_t = U_q \mathbf{w}_t / \|\mathbf{w}_t\|_q$; otherwise $\mathbf{w}'_t = \mathbf{w}_t$.

Thus, we swapped the notation from the explicit update and use \mathbf{w}'_t for the bounded and \mathbf{w}_t for the unbounded weight. Basically we now want to predict with the unbounded weights. The bound is as expected.

Theorem 5: Fix p and q such that $1/p + 1/q = 1$ and $2 \leq p < \infty$. Assume $\|\mathbf{x}_t\|_p \leq X_p$ and $\|\mathbf{u}_t\|_q \leq U_q$ for all t . Then the

bounded implicit update for Δ_q with learning rate $\eta = 1/((p-1)X_p^2)$ and parameter U_q satisfies

$$\begin{aligned} \sum_{t=1}^T (\mathbf{u}_t \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t)^2 &\leq \sum_{t=1}^T (\mathbf{u}_t \cdot \mathbf{x}_t - y_t)^2 \\ &+ (p-1)X_p^2 U_q^2 + 2(p-1)X_p^2 U_q \sum_{t=1}^{T-1} \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q. \end{aligned}$$

Proof: We mimic the proof of Theorem 4. This time we set

$$\begin{aligned} d_t^1 &= \Delta_q(\mathbf{u}_t, \mathbf{w}'_{t-1}) - \Delta_q(\mathbf{u}_t, \mathbf{w}_t) \\ d_t^2 &= \Delta_q(\mathbf{u}_t, \mathbf{w}_t) - \Delta_q(\mathbf{u}_t, \mathbf{w}'_t) \\ d_t^3 &= \Delta_q(\mathbf{u}_t, \mathbf{w}'_t) - \Delta_q(\mathbf{u}_{t+1}, \mathbf{w}'_t). \end{aligned}$$

For $\eta \leq 1/((p-1)X_p^2)$, the proof of Theorem 3 gives

$$d_t^1 \geq \frac{\eta}{2} s_t^2 - \frac{\eta}{2} r_t^2$$

where $s_t = \mathbf{u}_t \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t$ and $r_t = \mathbf{u}_t \cdot \mathbf{x}_t - y_t$. We estimate d_t^2 and d_t^3 and sum over t exactly as in the proof of Theorem 4. ■

All the previous bounds are for algorithms that use a constant learning rate that needs to be set at the beginning, and the optimal choice depends on the norms of the instances, which may not be known in advance. We close this section by considering a variant where we use a variable learning rate based on the norms of instances seen thus far. For simplicity, we deal only with the explicit update case.

Thus, define the *explicit update with variable learning rate* as

$$\mathbf{w}'_t = \mathbf{f}^{-1}(\mathbf{f}(\mathbf{w}_{t-1}) - \eta_t(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - y_t)\mathbf{x}_t)$$

where now η_t is a time-dependent learning rate. The out-of-bound update is as before.

The bound proven below is identical to the fixed η version given in Theorem 4 except for an additional factor of five in the second term on the right-hand side.

Theorem 6: Fix p and q such that $1/p + 1/q = 1$ and $2 \leq p < \infty$. Let $\eta_t = 1/((p-1)X_{p,t}^2)$ where $X_{p,t} = \max_{\tau \leq t} \|\mathbf{x}_\tau\|_p$. Assume $\|\mathbf{u}_t\|_q \leq U_q$ for all t . Then the bounded explicit update for Δ_q with the variable learning rate $\eta_t = 1/((p-1)X_{p,t}^2)$ and parameter U_q satisfies

$$\begin{aligned} \sum_{t=1}^T (\mathbf{u}_t \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 &\leq \sum_{t=1}^T (\mathbf{u}_t \cdot \mathbf{x}_t - y_t)^2 + 5(p-1)X_{p,T}^2 U_q^2 \\ &+ 2(p-1)X_{p,T}^2 U_q \sum_{t=1}^{T-1} \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q. \end{aligned}$$

Proof: We modify the proof of Theorem 4 using the method of [18] for handling the variable learning rate. Fortunately, in filtering, the technicalities are much easier than in the prediction setting.

Thus, we consider the quantity $\Delta_q(\mathbf{u}_t, \mathbf{w}_{t-1})/\eta_t$. By replacing η with η_t in (20), we see that the proof of Theorem 4 implies

$$\begin{aligned} \frac{\Delta_q(\mathbf{u}_t, \mathbf{w}_{t-1})}{\eta_t} - \frac{\Delta_q(\mathbf{u}_{t+1}, \mathbf{w}_t)}{\eta_t} &\geq \frac{s_t^2}{2} - \frac{r_t^2}{2} + \frac{1}{2\eta_t} \|\mathbf{u}_t\|_q^2 \\ &- \frac{1}{2\eta_t} \|\mathbf{u}_{t+1}\|_q^2 - \frac{U_q}{\eta_t} \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q \end{aligned}$$

where $s_t = \mathbf{u}_t \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t$ and $r_t = \mathbf{u}_t \cdot \mathbf{x}_t - y_t$. (Again we set $\mathbf{u}_{T+1} = \mathbf{u}_T$; also let $X_{p,T+1} = X_{p,T}$.) By substituting $\eta_t = 1/((p-1)X_{p,t}^2)$ and then noticing that $X_{p,t} \leq X_{p,t+1} \leq X_{p,T}$, we get

$$\begin{aligned} &(p-1)X_{p,t}^2 (\Delta_q(\mathbf{u}_t, \mathbf{w}_{t-1}) - \Delta_q(\mathbf{u}_{t+1}, \mathbf{w}_t)) \\ &\geq \frac{s_t^2}{2} - \frac{r_t^2}{2} + \frac{p-1}{2} X_{p,t}^2 \|\mathbf{u}_t\|_q^2 - \frac{p-1}{2} X_{p,t+1}^2 \|\mathbf{u}_{t+1}\|_q^2 \\ &- (p-1)X_{p,T}^2 U_q \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q. \end{aligned}$$

By [18, Lemma 3.2], we have $\Delta_q(\mathbf{v}, \mathbf{v}') \leq 2V^2$ whenever $\|\mathbf{v}\|_q \leq V$ and $\|\mathbf{v}'\|_q \leq V$, so in particular $\Delta_q(\mathbf{u}_{t+1}, \mathbf{w}_t) \leq 2U_q^2$. Remembering that $X_{p,t}^2 - X_{p,t+1}^2 \leq 0$, we get

$$\begin{aligned} &(p-1)X_{p,t}^2 \Delta_q(\mathbf{u}_t, \mathbf{w}_{t-1}) - (p-1)X_{p,t+1}^2 \Delta_q(\mathbf{u}_{t+1}, \mathbf{w}_t) \\ &= (p-1)X_{p,t}^2 (\Delta_q(\mathbf{u}_t, \mathbf{w}_{t-1}) - \Delta_q(\mathbf{u}_{t+1}, \mathbf{w}_t)) \\ &+ (p-1)(X_{p,t}^2 - X_{p,t+1}^2) \Delta_q(\mathbf{u}_{t+1}, \mathbf{w}_t) \\ &\geq \frac{s_t^2}{2} - \frac{r_t^2}{2} + \frac{p-1}{2} X_{p,t}^2 \|\mathbf{u}_t\|_q^2 - \frac{p-1}{2} X_{p,t+1}^2 \|\mathbf{u}_{t+1}\|_q^2 \\ &- (p-1)X_{p,T}^2 U_q \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q \\ &+ 2(p-1)(X_{p,t}^2 - X_{p,t+1}^2) U_q^2. \end{aligned}$$

By summing over $t = 1, \dots, T$, we get

$$\begin{aligned} &(p-1)X_{p,1}^2 \Delta_q(\mathbf{u}_1, \mathbf{w}_0) - (p-1)X_{p,T+1}^2 \Delta_q(\mathbf{u}_{T+1}, \mathbf{w}_T) \\ &\geq \frac{1}{2} \sum_{t=1}^T s_t^2 - \frac{1}{2} \sum_{t=1}^T r_t^2 + \frac{p-1}{2} X_{p,1}^2 \|\mathbf{u}_1\|_q^2 \\ &- \frac{p-1}{2} X_{p,T+1}^2 \|\mathbf{u}_{T+1}\|_q^2 - (p-1)X_{p,T}^2 U_q \\ &\times \sum_{t=1}^T \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q + 2(p-1)(X_{p,1}^2 - X_{p,T+1}^2) U_q^2. \end{aligned}$$

The result follows by solving for $\sum_{t=1}^T s_t^2$, noticing $\Delta_q(\mathbf{u}_1, \mathbf{w}_0) = \|\mathbf{u}_1\|_q^2/2$ and then ignoring the negative terms $-(p-1)X_{p,T+1}^2 \Delta_q(\mathbf{u}_{T+1}, \mathbf{w}_T)$ and $4(p-1)X_{p,1}^2 U_q^2$. ■

VI. GENERALIZED LINEAR MODELS

We extended framework slightly to cover *generalized linear regression*. Here we replace the model (1) by

$$y_t = h(\mathbf{u} \cdot \mathbf{x}_t + v_t) \quad (21)$$

where h is a continuous, strictly increasing *transfer function*. The logistic sigmoid $h(r) = 1/(1 + \exp(-r))$ is a typical example. In the prediction setting (where the learner tries to match y_t), the prediction becomes $\hat{y}_t = h(\mathbf{w}_{t-1} \cdot \mathbf{x}_t)$. In the filtering setting, we would naturally also include the transfer function in the prediction, giving $\hat{y}_t = h(\mathbf{w}_{t-1} \cdot \mathbf{x}_t)$ for the *a priori* and $\hat{y}_t = h(\mathbf{w}_t \cdot \mathbf{x}_t)$ for the *a posteriori* case. The algorithm then tries to match \hat{y}_t to $h(\mathbf{u} \cdot \mathbf{x}_t)$. One could in principle still use the squared error $(h(\mathbf{u} \cdot \mathbf{x}_t) - h(\mathbf{w} \cdot \mathbf{x}_t))^2$ as the performance measure, but this is nonconvex in \mathbf{u} and \mathbf{w} and actually leads to a very badly behaved optimization problem [19]. We obtain a better behaved problem by using the *matching loss* for h [19], defined for y and y' in the range of h as

$$L(y, y') = \int_{h^{-1}(y)}^{h^{-1}(y')} (h(r) - y) dr. \quad (22)$$

(Notice that by our assumptions h is one-to-one.) It is easy to see that for the identity transfer function $h(r) = r$, we get

$L(y, y') = (y - y')^2/2$; and for the logistic sigmoid $h(r) = 1/(1 + \exp(-r))$, we get the logarithmic loss

$$L(y, y') = y \ln \frac{y}{y'} + (1 - y) \ln \frac{1 - y}{1 - y'}.$$

The definition (22) may seem arbitrary, but it is actually a one-dimensional Bregman divergence: if we let $H(r) = \int h(r) dr$, then

$$L(h(a), h(a')) = \Delta_H(a', a). \quad (23)$$

Using a Bregman divergence as a loss naturally generalizes to multidimensional outputs [6], but we shall not pursue that here.

Directly from (22), we obtain a simple expression for its gradient

$$\nabla_{\mathbf{w}} L(y, h(\mathbf{w} \cdot \mathbf{x})) = (h(\mathbf{w} \cdot \mathbf{x}) - y) \mathbf{x}. \quad (24)$$

Therefore, the explicit update (14) naturally generalizes to

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

where $\hat{y}_t = h(\mathbf{w}_{t-1} \cdot \mathbf{x}_t)$. The implicit update can be generalized similarly; for it we use $\hat{y}_t = h(\mathbf{w}_t \cdot \mathbf{x}_t)$. For these updates we can now prove bounds that have as an additional factor an upper bound on the slope of the transfer function. The techniques are essentially those introduced by [10].

Theorem 7: Fix p and q such that $1/p + 1/q = 1$ and $2 \leq p < \infty$. Let h be strictly increasing and continuously differentiable with c such that $0 < h(r) \leq c$ holds for all r , and let L be the matching loss for h . Assume that $\|\mathbf{x}_t\|_p \leq X_p$ for all t . Then both the explicit update and implicit update for Δ_q with learning rate $\eta = 1/((p-1)cX_p^2)$ satisfy

$$\sum_{t=1}^T L(\hat{y}_t, h(\mathbf{u} \cdot \mathbf{x}_t)) \leq \sum_{t=1}^T L(y_t, h(\mathbf{u} \cdot \mathbf{x}_t)) + (p-1)cX_p^2 \|\mathbf{u}\|_q^2$$

for any $\mathbf{u} \in \mathbf{R}^n$.

Proof: Consider first the explicit update. As in the proof of Theorem 2, let

$$\begin{aligned} d_t &= \Delta_q(\mathbf{u}, \mathbf{w}_{t-1}) - \Delta_q(\mathbf{u}, \mathbf{w}_t) \\ &= \eta(y_t - \hat{y}_t)(\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t) - \Delta_q(\mathbf{w}_{t-1}, \mathbf{w}_t). \end{aligned}$$

Using (23), we get

$$\begin{aligned} (y_t - \hat{y}_t)(\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t) &= L(\hat{y}_t, h(\mathbf{u} \cdot \mathbf{x}_t)) \\ &\quad - L(y_t, h(\mathbf{u} \cdot \mathbf{x}_t)) + L(y_t, \hat{y}_t). \end{aligned}$$

Simple calculus shows that $L(y, \hat{y}) \geq (1/2c)(y - \hat{y})^2$ for all y and \hat{y} . By combining this with (15), we get

$$\begin{aligned} d_t &\geq \eta(L(\hat{y}_t, h(\mathbf{u} \cdot \mathbf{x}_t)) - L(y_t, h(\mathbf{u} \cdot \mathbf{x}_t))) \\ &\quad + \frac{\eta}{2}(y_t - \hat{y}_t)^2 \left(\frac{1}{c} - \eta(p-1)X_p^2 \right). \end{aligned}$$

The claim follows by summing over t as usual.

Consider now the implicit update. We have

$$\begin{aligned} d_t &= \Delta_q(\mathbf{u}, \mathbf{w}_{t-1}) - \Delta_q(\mathbf{u}, \mathbf{w}_t) \\ &= \eta(y_t - \hat{y}_t)(\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t) - \Delta_q(\mathbf{w}_{t-1}, \mathbf{w}_t) \end{aligned}$$

where now $\hat{y}_t = h(\mathbf{w}_t \cdot \mathbf{x}_t)$. We write

$$\begin{aligned} (y_t - \hat{y}_t)(\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t) &= (y_t - \hat{y}_t)(\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t) \\ &\quad + (y_t - \hat{y}_t)(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t). \end{aligned}$$

Like above, we have

$$(y_t - \hat{y}_t)(\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t) = L(\hat{y}_t, h(\mathbf{u} \cdot \mathbf{x}_t))$$

$$-L(y_t, h(\mathbf{u} \cdot \mathbf{x}_t)) + L(y_t, \hat{y}_t).$$

Also, since \mathbf{w}_t is the solution to

$$\mathbf{w}_t = \arg \min_{\mathbf{w}} (\Delta_q(\mathbf{w}, \mathbf{w}_{t-1}) + \eta L(y_t, h(\mathbf{w} \cdot \mathbf{x}_t)))$$

we have either $y_t \leq \mathbf{w}_t \cdot \mathbf{x}_t \leq \mathbf{w}_{t-1} \cdot \mathbf{x}_t$ or $\mathbf{w}_{t-1} \cdot \mathbf{x}_t \leq \mathbf{w}_t \cdot \mathbf{x}_t \leq y_t$. In either case, $(y_t - \hat{y}_t)(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - \mathbf{w}_t \cdot \mathbf{x}_t) \leq 0$. Hence, we have established

$$\begin{aligned} d_t &\geq \eta(L(\hat{y}_t, h(\mathbf{u} \cdot \mathbf{x}_t)) - L(y_t, h(\mathbf{u} \cdot \mathbf{x}_t)) \\ &\quad + L(y_t, \hat{y}_t)) - \Delta_q(\mathbf{w}_{t-1}, \mathbf{w}_t) \end{aligned}$$

and can proceed as with the explicit update. \blacksquare

Because of how we defined \hat{y}_t , the theorem gives an *a priori* filtering bound for the explicit update and an *a posteriori* bound for the implicit update.

When h is the identity function, we get the results of Section IV with $c = 1$. For the logistic sigmoid, $c = 1/4$. Thresholded transfer functions, such as $h(r) = \text{sign}(r)$, correspond to the limiting case $c \rightarrow \infty$, which makes the bound vacuous.

This result generalizes to the nonstationary case (Section V) in the obvious manner; we omit the details.

Our main motivation for considering loss functions other than square loss was that they make the problem involving a nonlinear transfer function computationally simpler, which also allows strong worst case bounds. One might also prefer different loss functions if one assumes a non-Gaussian noise distribution [20]. This is quite different from our framework, where no statistical assumptions are made.

VII. SIMULATION RESULTS

The discussion following Theorem 2 suggests that having a sparse target favors having a large p . We illustrate this with a simple filtering simulation.

At time t , the sender sends a bit $y_t \in \{-1, 1\}$ over a channel. The recipient is required to produce a binary prediction $\hat{y}_t \in \{-1, 1\}$ about the sent bit. If $\hat{y}_t \neq y_t$, we say that an error occurred. What the recipient actually observes is

$$r_t = \sum_{i=0}^{k-1} u_{i+1} y_{t-i} + v_t$$

where $\mathbf{u} \in \mathbf{R}^k$ for some k describes the channel and v_t is zero-mean Gaussian noise. The prediction is then $\hat{y}_t = \text{sign}(\mathbf{w}_{t-1} \cdot \mathbf{x}_t)$, where $\mathbf{x}_t = (r_{t-m}, \dots, r_t, \dots, r_{t+m}) \in \mathbf{R}^{2m+1}$ and $n = 2m + 1$ is the filter length.

Notice that this setting is not quite the same as introduced earlier, since we are now considering discrete errors but still using the update rules based on square loss. The purpose of this is to illustrate how the algorithms work on binary prediction, which often is the problem one is really interested in.

For choosing \mathbf{u} , we considered two different distributions. In the first experiment, \mathbf{u} is from a Gaussian with unit variance. In the second experiment, $u_i = s_i e^{r_i}$, where $s_i \in \{-1, 1\}$ and $r_i \in [-10, 10]$ are distributed uniformly. In both cases, we then renormalize to make $\|\mathbf{u}\|_2 = 1$. The targets \mathbf{u} from the second distribution are ‘‘sparse’’ in the sense that most of the weight is concentrated on only few components, whereas the targets from the first distribution are ‘‘dense.’’ In both experiments, we

used $k = 10$, $m = 15$ and a signal-to-noise ratio of 10 dB. We compared the explicit update algorithm with $p = 2$ against $p = 2 \ln n \approx 6.9$. (As we remarked after Theorem 2, for $p = 2 \ln n$ we can estimate $(p - 1)\|\mathbf{x}\|_p^2 \leq 2e(\ln n)\|\mathbf{x}\|_\infty^2$.)

Notice that due to the constant learning rate, the weight vectors of the algorithms end up oscillating around the optimum, so the algorithms converge to a nonzero error rate. By using a smaller learning rate, one can reduce the oscillations and thus achieve a smaller final error rate, but this makes the initial convergence slower. The choice of learning rate is thus not straightforward.

We used for $p = 2 \ln n$ the value $\eta = 1/((p - 1)X_p^2)$ as suggested by Theorem 2. This gave final error rates 0.02 in the first experiment and 0.01 in the second one. For $p = 2$ we then chose η so that these same final error rates were achieved. For the first experiment, this resulted in $\eta = 0.45/X_2^2$, and for the second one, $\eta = 0.4/X_2^2$.

The development of the error rates over time is shown in Fig. 1. As expected, $p = 2$ gives a faster convergence for dense targets and $p = 2 \ln n$ for sparse targets. The differences here are not large, but they become more apparent if the filter length (i.e., dimensionality of inputs) is increased.

We did not include the implicit updates in this comparison. In other experiments we noticed that for any fixed p and η , the implicit update has slower initial convergence and smaller final error rate than the explicit one. This can be understood by noticing that by (17), the implicit update always makes a smaller step. Hence, as a crude first approximation, the implicit update is similar to the explicit update with a smaller η .

VIII. DISCUSSION AND CONCLUSION

We have shown how Bregman divergences based on p -norms can be used to derive generalizations of the classical LMS algorithm. This is a direct application of methods recently introduced in machine learning. The resulting p -norm algorithms have for large p quite different behavior from the LMS, which is the special case $p = 2$. In particular, both theoretical bounds and preliminary simulations suggest that the large p version has better performance when the target weight vector is sparse. We apply further methods from machine learning to show that also in filtering, the p -norm algorithms can be made robust against target shift and can be adapted for generalized linear systems.

The question of applying these techniques to genuinely nonlinear problems remains unsolved. Recently much work has been done in machine learning on applying linear algorithm to nonlinear problems using the so-called kernel trick. This trick works for a large class of algorithms, such as LMS, the support vector machine, or more generally any *rotation invariant* algorithm [16], [17], [21]. The p -norm algorithm for $p \neq 2$ is not rotation invariant, and it remains an open problem whether it can be efficiently nonlinearized with some technique analogous to the kernel trick. For algorithms with similar performance to the p -norm algorithm with large p , efficient techniques have

been found for some kernels [22], but for other kernels the problem is known to be intractable [23]. Further, the computational requirements in signal processing applications may even rule out kernel-style approaches that rely on storing a large number of data points. Thus, the prospects of finding a general nonlinear version of the p -norm algorithms do not seem good.

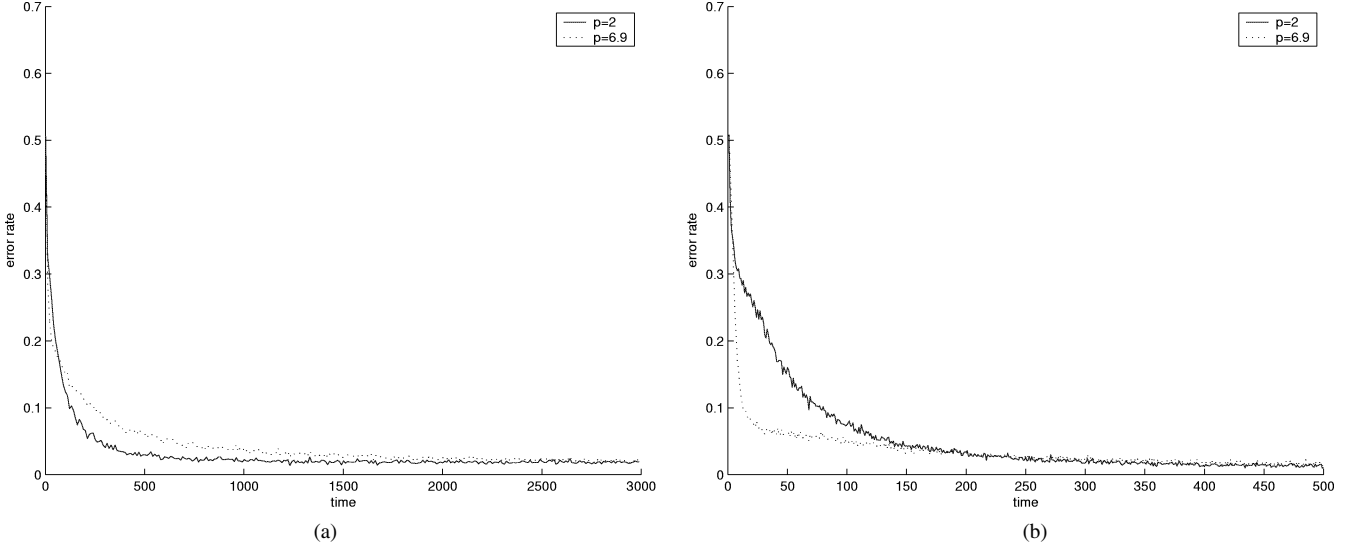


Fig. 1. Error rates as function of time for $p = 2$ (solid line) and $p \approx 6.9$ (dotted line) in the filtering simulation. The error rates are averages over 5000 runs. The experiments were run for 30000 time steps to make sure the algorithms converged to same error rate; the plots show only the initial part. (a) Dense target and (b) sparse target.

APPENDIX I PROOF OF (15)

Since $1/p + 1/q = 1$, a straightforward calculation shows that $\|\mathbf{w}\|_p = \|\mathbf{f}(\mathbf{w})\|_q$ and $\mathbf{w} \cdot \mathbf{f}(\mathbf{w}) = \|\mathbf{w}\|_q^2$ for all $\mathbf{w} \in \mathbf{R}^n$ [8, Lemma 1]. Fix now $\boldsymbol{\theta} = \mathbf{f}(\mathbf{w})$ and $\boldsymbol{\theta}' = \mathbf{f}(\mathbf{w}')$, with $\mathbf{x} = \boldsymbol{\theta}' - \boldsymbol{\theta}$. Based on the above, it is easy to verify that $\Delta_q(\mathbf{w}, \mathbf{w}') = \Delta_p(\boldsymbol{\theta}', \boldsymbol{\theta})$. (Notice the order of the arguments.) Let $G(\boldsymbol{\theta}) = (1/2)\|\boldsymbol{\theta}\|_p^2$. Since Δ_p is defined as the error of a first-order Taylor approximation for G , we can write

$$\Delta_p(\boldsymbol{\theta} + \mathbf{x}, \boldsymbol{\theta}) = \frac{1}{2} \mathbf{x}^T H \mathbf{x} \quad (25)$$

where $H_{ij} = \partial^2 G(\boldsymbol{\xi}) / \partial \xi_i \partial \xi_j$ and the derivatives are evaluated at some point $\boldsymbol{\xi}$ on the line between $\boldsymbol{\theta}$ and $\boldsymbol{\theta} + \mathbf{x}$. We now estimate the right-hand side of (25) as in [7, Theorem 7.1]. We have

$$H_{ij} = (2-p) \text{sign}(\xi_i) |\xi_i|^{p-1} \text{sign}(\xi_j) |\xi_j|^{p-1} \|\boldsymbol{\xi}\|_p^{2-2p} + \delta_{ij} (p-1) |\xi_i|^{p-2} \|\boldsymbol{\xi}\|_p^{2-p}.$$

Since we assume $p \geq 2$, we get

$$\begin{aligned} \mathbf{x}^T H \mathbf{x} &= (2-p) \|\boldsymbol{\xi}\|_p^{2-2p} \left(\sum_i \text{sign}(\xi_i) |\xi_i|^{p-1} x_i \right)^2 \\ &\quad + (p-1) \|\boldsymbol{\xi}\|_p^{2-p} \left(\sum_i |\xi_i|^{p-2} x_i^2 \right) \\ &\leq (p-1) \|\boldsymbol{\xi}\|_p^{2-p} \tilde{\boldsymbol{\xi}} \cdot \tilde{\mathbf{x}} \end{aligned}$$

where $\tilde{\xi}_i = |\xi_i|^{p-2}$ and $\tilde{x}_i = x_i^2$. Since $(1/(p/(p-2))) + (1/(p/2)) = 1$, Hölder's inequality gives us

$$|\tilde{\boldsymbol{\xi}} \cdot \tilde{\mathbf{x}}| \leq \|\tilde{\boldsymbol{\xi}}\|_{\frac{p}{p-2}} \|\tilde{\mathbf{x}}\|_{\frac{p}{2}} = \|\boldsymbol{\xi}\|_p^{p-2} \|\mathbf{x}\|_p^2$$

and the claim follows.

APPENDIX II EXPONENTIATED GRADIENT

As in Example 2, the relative entropy can be seen as a Bregman divergence. The constraint $\sum_i w_i = 1$ requires some additional technicalities. We present here a fairly straightforward method. For a more general framework allowing potential functions F that are not *strictly* convex, see [6].

For \mathbf{w} and \mathbf{w}' with $w_i, w'_i \geq 0$, and $\sum_i w_i = \sum_i w'_i = 1$, define the relative entropy

$$\Delta_{\text{re}}(\mathbf{w}, \mathbf{w}') = \sum_{i=1}^n w_i \ln \frac{w_i}{w'_i}.$$

(We take $0 \ln 0 = 0$ and $\ln 0 = \infty$ otherwise.) Notice that $\Delta_{\text{re}}(\mathbf{w}, \mathbf{w}')$ is convex in \mathbf{w} . Consider minimizing

$$C_t(\mathbf{w}) = \Delta_{\text{re}}(\mathbf{w}, \mathbf{w}_{t-1}) + \frac{1}{2} \eta (y_t - \mathbf{w} \cdot \mathbf{x}_t)^2$$

subject to $\sum_i w_i = 1$, $w_i \geq 0$ for all i . The problem is convex, so we solve it by setting the gradient of the Lagrangian

$$A_t(\mathbf{w}) = \Delta_{\text{re}}(\mathbf{w}, \mathbf{w}_{t-1}) + \lambda \left(\sum_{i=1}^n w_i - 1 \right) + \frac{1}{2} \eta (y_t - \mathbf{w} \cdot \mathbf{x}_t)^2$$

to zero. This yields

$$\ln \frac{w_i}{w_{t-1,i}} + 1 + \lambda + \eta (\mathbf{w} \cdot \mathbf{x}_t - y_t) x_{t,i} = 0$$

or (after substituting λ such that $\sum_i w_i = 1$)

$$w_i = \frac{w_{t-1,i} \exp(-\eta (\mathbf{w} \cdot \mathbf{x}_t - y_t) x_{t,i})}{Z} \quad (26)$$

where $Z = \sum_{i=1}^n w_{t-1,i} \exp(-\eta (\mathbf{w} \cdot \mathbf{x}_t - y_t) x_{t,i})$. Notice that $w_{t,i} > 0$ implies $w_i > 0$.

For $\mathbf{z} \in \mathbf{R}^n$, define now $\mathbf{g}(\mathbf{z})$ by

$$g_i(\mathbf{z}) = \frac{e^{z_i}}{\sum_{j=1}^n e^{z_j}}. \quad (27)$$

Let \mathbf{z}_{t-1} be such that $\mathbf{g}(\mathbf{z}_{t-1}) = \mathbf{w}_{t-1}$. It is easy to see that such a \mathbf{z}_{t-1} exists assuming $w_{t-1,i} > 0$ for all i and $\sum_{i=1}^n w_{t-1,i} = 1$. Further, if \mathbf{z}'_{t-1} is another vector satisfying $\mathbf{g}(\mathbf{z}'_{t-1}) = \mathbf{w}_{t-1}$, then \mathbf{z}_{t-1} and \mathbf{z}'_{t-1} are the same up to an additive constant, i.e., $z_{t-1,i} = z'_{t-1,i} + b$ for some b that does not depend on i . Equation (26) can now be written as $\mathbf{w} = \mathbf{g}(\mathbf{z})$, where $\mathbf{z} = \mathbf{z}_{t-1} - \eta(\mathbf{w} \cdot \mathbf{x}_t - y_t)\mathbf{x}_t$. Notice that because of the normalization, the choice of the representative \mathbf{z}_{t-1} (i.e., the constant b) makes no difference.

Again, we define the implicit and explicit version of the update. We use an additional parameter vector \mathbf{z}_t to present the algorithm, the actual weights being given by $\mathbf{w}_t = \mathbf{g}(\mathbf{z}_t)$. In both cases, we start with $\mathbf{z}_1 = \mathbf{0}$. For the *implicit exponentiated gradient* algorithm, we define \mathbf{z}_t by

$$\mathbf{z}_t = \mathbf{z}_{t-1} - \eta(\mathbf{w}_t \cdot \mathbf{x}_t - y_t)\mathbf{x}_t$$

and for *explicit exponentiated gradient* (EG) algorithm by

$$\mathbf{z}_t = \mathbf{z}_{t-1} - \eta(\mathbf{w}_{t-1} \cdot \mathbf{x}_t - y_t)\mathbf{x}_t.$$

Thus the implicit update uses as \mathbf{w}_t the minimizer of C_t , while the explicit update uses an approximation thereof. These updates are analogous to the implicit and explicit updates given previously, with \mathbf{g} now replacing \mathbf{f}^{-1} . However, in this case \mathbf{g} is not one-to-one, so we write the update in terms of \mathbf{z}_t (which corresponds to $\mathbf{f}(\mathbf{w}_t)$ in the previous setting) and not directly in terms of \mathbf{w}_t .

The following lemma gives the analogues of (11) and (15) for relative entropy.

Lemma 1: Let $\mathbf{w} = \mathbf{g}(\mathbf{z})$ and $\mathbf{w}' = \mathbf{g}(\mathbf{z}')$ for some $\mathbf{z}, \mathbf{z}' \in \mathbf{R}^n$. Then

$$\Delta_{\text{re}}(\mathbf{w}, \mathbf{w}') \leq \frac{1}{8} \left(\max_i (z'_i - z_i) - \min_i (z'_i - z_i) \right)^2 \quad (28)$$

and for any $\mathbf{u} \in \mathbf{R}^n$ with $u_i \geq 0$ and $\sum_i u_i = 1$, we have

$$\Delta_{\text{re}}(\mathbf{u}, \mathbf{w}') = \Delta_{\text{re}}(\mathbf{u}, \mathbf{w}) + \Delta_{\text{re}}(\mathbf{w}, \mathbf{w}') + (\mathbf{z}' - \mathbf{z}) \cdot (\mathbf{w} - \mathbf{u}). \quad (29)$$

Proof: Equation (29) follows directly from the definition. To prove (28), we first write

$$\Delta_{\text{re}}(\mathbf{w}, \mathbf{w}') = G(\mathbf{z}') - (G(\mathbf{z}) + \mathbf{g}(\mathbf{z}) \cdot (\mathbf{z}' - \mathbf{z}))$$

where $G(\mathbf{z}) = \ln(\sum_i e^{z_i})$. Notice that $\mathbf{g} = \nabla G$. Therefore, $\Delta_{\text{re}}(\mathbf{w}, \mathbf{w}')$ is the error in the first-order Taylor approximation of $G(\mathbf{z}')$ around $G(\mathbf{z})$, and we have $\Delta_{\text{re}}(\mathbf{w}, \mathbf{w}') = (1/2)(\mathbf{z}' - \mathbf{z})^T H(\mathbf{z}' - \mathbf{z})$, where H is the Hessian of G evaluated at some point between \mathbf{z} and \mathbf{z}' . We have

$$\frac{\partial^2 G(\mathbf{z})}{\partial z_i \partial z_j} = \frac{\partial g_i(\mathbf{z})}{\partial z_j} = \delta_{ij} g_i(\mathbf{z}) - g_i(\mathbf{z}) g_j(\mathbf{z}).$$

Therefore we can write $H_{ij} = \delta_{ij} p_i - p_i p_j$ for some \mathbf{p} that satisfies $p_i > 0$ and $\sum_i p_i = 1$. Denote now by X a random

variable that is obtained by choosing the value $x_i = z'_i - z_i$ with probability p_i . Then

$$\begin{aligned} (\mathbf{z}' - \mathbf{z})^T H(\mathbf{z}' - \mathbf{z}) &= \sum_i p_i x_i^2 - \sum_{i,j} x_i x_j p_i p_j \\ &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 \\ &= \text{Var}[X] \\ &\leq \frac{1}{4} (\max_i x_i - \min_i x_i)^2. \end{aligned}$$

Theorem 8: Assume that $\max_i x_{t,i} - \min_i x_{t,i} \leq R$ for all t . Then for any $\mathbf{u} \in \mathbf{R}^n$ with $u_i \geq 0$ and $\sum_i u_i = 1$, the explicit EG algorithm with learning rate $\eta = 4/R^2$ satisfies

$$\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 \leq \sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + \frac{1}{4} R^2 \Delta_{\text{re}}(\mathbf{u}, \mathbf{w}_0)$$

where $\mathbf{w}_0 = \mathbf{g}(\mathbf{0})$ is the uniform weight vector.

Proof: We analyze the progress $d_t = \Delta_{\text{re}}(\mathbf{u}, \mathbf{w}_{t-1}) - \Delta_{\text{re}}(\mathbf{u}, \mathbf{w}_t)$. By substituting the explicit EG update into (29) and then using (28), we get

$$\begin{aligned} d_t &= \eta(y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)\mathbf{x}_t \cdot (\mathbf{u}_t - \mathbf{w}_{t-1}) - \Delta_{\text{re}}(\mathbf{w}_{t-1}, \mathbf{w}_t) \\ &\geq \eta(y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)(\mathbf{u}_t \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t) \\ &\quad - \frac{1}{8} \eta^2 (y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 R^2. \end{aligned}$$

By rearranging terms, we can write this as

$$d_t \geq \frac{\eta}{2} s_t^2 - \frac{\eta}{2} r_t^2 + \frac{\eta}{2} (s_t - r_t)^2 \left(1 - \frac{\eta R^2}{4} \right)$$

where $s_t = \mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t$ and $r_t = \mathbf{u} \cdot \mathbf{x}_t - y_t$. Since $\eta R^2/4 = 1$, we can apply $\Delta_{\text{re}}(\mathbf{u}, \mathbf{w}_{T+1}) \geq 0$ to get

$$\begin{aligned} \Delta_{\text{re}}(\mathbf{u}, \mathbf{w}_0) &\geq \Delta_{\text{re}}(\mathbf{u}, \mathbf{w}_0) - \Delta_{\text{re}}(\mathbf{u}, \mathbf{w}_{T+1}) \\ &= \sum_{t=1}^T d_t \\ &\geq \frac{4}{R^2} \left(\sum_{t=1}^T s_t^2 - \sum_{t=1}^T r_t^2 \right) \end{aligned}$$

from which the claim follows. \blacksquare

The above theorem assumes the comparison vector \mathbf{u} is a probability vector. To deal with arbitrary vectors \mathbf{u} with $\|\mathbf{u}\|_1 \leq U_1$ for some given bound $U_1 > 0$, we define the *scaled explicit EG $^\pm$* algorithm as explicit EG with each input \mathbf{x}_t replaced by $\mathbf{x}'_t = (Ux_{t,1}, \dots, Ux_{t,n}, -Ux_{t,1}, \dots, -Ux_{t,n}) \in \mathbf{R}^{2n}$.

Corollary 1: Assume $\|\mathbf{x}_t\|_\infty \leq X_\infty$ for all t . Then for any $\mathbf{u} \in \mathbf{R}^n$ with $\|\mathbf{u}\|_1 \leq U_1$, the scaled explicit EG $^\pm$ algorithm satisfies

$$\sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 \leq \sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - y_t)^2 + \ln(2n) X_\infty^2 U_1^2.$$

Proof: There is some $\mathbf{u}' \in \mathbf{R}^{2n}$ with $u'_i \geq 0$ for all i and $\sum_i u'_i = 1$ such that $\mathbf{u}' \cdot \mathbf{x}'_t = \mathbf{u} \cdot \mathbf{x}_t$ for all t . Thus we can apply Theorem 8 with this \mathbf{u}' . We have $\max_i x'_{t,i} - \min_i x'_{t,i} = 2U_1 \|\mathbf{x}_t\|_\infty$. Since \mathbf{w}_0 is the uniform $2n$ -dimensional probability vector, we have $\Delta_{\text{re}}(\mathbf{u}', \mathbf{w}_0) \leq \ln(2n)$. \blacksquare

Bounds for the implicit EG algorithm can be proven analogously.

REFERENCES

- [1] J. Kivinen, M. K. Warmuth, and B. Hassibi, "The p -norm generalization of the LMS algorithm for adaptive filtering," presented at the 13th IFAC Symp. System Identification, P. M. J. V. den Hof, B. Wahlberg, and S. Weiland, Eds., Rotterdam, The Netherlands, Aug. 27–29, 2003.
- [2] B. Widrow and M. E. Hoff, "Adaptive switching circuits," in *1960 IRE WESCON Convention Rec.*, 1960, pp. 96–104.
- [3] B. Hassibi, A. H. Sayed, and T. Kailath, " H^∞ optimality of the LMS algorithm," *IEEE Trans. Signal Process.*, vol. 44, no. 2, pp. 267–280, Feb. 1996.
- [4] N. Cesa-Bianchi, P. Long, and M. Warmuth, "Worst-case quadratic loss bounds for on-line prediction of linear functions by gradient descent," *IEEE Trans. Neural Netw.*, vol. 7, no. 3, pp. 604–619, May 1996.
- [5] J. Kivinen and M. K. Warmuth, "Additive versus exponentiated gradient updates for linear prediction," *Inf. Comput.*, vol. 132, no. 1, pp. 1–64, Jan. 1997.
- [6] —, "Relative loss bounds for multidimensional regression problems," *Mach. Learn.*, vol. 45, no. 3, pp. 301–329, Dec. 2001.
- [7] A. J. Grove, N. Littlestone, and D. Schuurmans, "General convergence results for linear discriminant updates," *Mach. Learn.*, vol. 43, no. 3, pp. 173–210, Jun. 2001.
- [8] C. Gentile and N. Littlestone, "The robustness of the p -norm algorithms," in *12th Annu. Conf. Comput. Learn. Theory*. New York: ACM Press, Jul. 1999, pp. 1–11.
- [9] M. Herbster and M. K. Warmuth, "Tracking the best linear predictor," *J. Mach. Learn. Res.*, vol. 1, pp. 281–309, Sep. 2001.
- [10] D. P. Helmbold, J. Kivinen, and M. K. Warmuth, "Relative loss bounds for single neurons," *IEEE Trans. Neural Netw.*, vol. 10, no. 6, pp. 1291–1304, Nov. 1999.
- [11] N. Cesa-Bianchi, A. Conconi, and C. Gentile, "On the generalization ability of on-line learning algorithms," *IEEE Trans. Inf. Theory*, vol. 50, no. 9, pp. 2050–2057, Sep. 2004.
- [12] K. S. Azoury and M. K. Warmuth, "Relative loss bounds for on-line density estimation with the exponential family of distributions," *Mach. Learn.*, vol. 43, no. 3, pp. 211–246, Jun. 2001.
- [13] L. Bregman, "The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming," *USSR Comp. Math. Phys.*, vol. 7, pp. 200–217, 1967.
- [14] R. E. Schapire and M. K. Warmuth, "On the worst-case analysis of temporal-difference learning algorithms," *Mach. Learn.*, vol. 22, no. 1/2/3, pp. 95–121, Jan. 1996.
- [15] N. Littlestone, "Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm," *Mach. Learn.*, vol. 2, no. 4, pp. 285–318, Apr. 1988.
- [16] J. Kivinen, M. K. Warmuth, and P. Auer, "The Perceptron algorithm vs. Winnow: Linear vs. logarithmic mistake bounds when few input variables are relevant," *Artif. Intell.*, vol. 97, pp. 325–343, Dec. 1997.
- [17] A. Y. Ng, "Feature selection, L_1 vs. L_2 regularization, and rotational invariance," in *Proc. 21st Int. Conf. Machine Learning*, R. Greiner and D. Schuurmans, Eds., Jul. 2004, pp. 615–622.
- [18] P. Auer, N. Cesa-Bianchi, and C. Gentile, "Adaptive and self-confident on-line learning algorithms," *J. Comp. Syst. Sci.*, vol. 64, no. 1, pp. 48–75, Feb. 2002.
- [19] P. Auer, M. Herbster, and M. K. Warmuth, "Exponentially many local minima for single neurons," in *Advances in Neural Information Processing Systems 8*, D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, Eds. Cambridge, MA: MIT Press, June 1996, pp. 316–317.
- [20] T. Y. Al-Naffouri, A. H. Sayed, and T. Kailath, "On the selection of optimal nonlinearities for stochastic gradient adaptive algorithms," in *Proc. 2000 IEEE Conf. Acoustics, Speech, Signal Processing (ICASSP)*, vol. 1, Jun. 2000, pp. 464–467.
- [21] M. K. Warmuth and S. V. N. Vishwanathan, "Leaving the span," in *Proc. 18th Annu. Conf. Learning Theory*, P. Auer and R. Meir, Eds. Berlin, Germany: Springer, Jun. 2005, pp. 366–381.
- [22] E. Takimoto and M. K. Warmuth, "Path kernels and multiplicative updates," *J. Mach. Learn. Res.*, vol. 4, pp. 773–818, Oct. 2003.
- [23] R. Khardon, D. Roth, and R. A. Servedio, "Efficiency versus convergence of boolean kernels for on-line learning algorithms," *J. Artif. Intell. Res.*, vol. 24, pp. 341–356, Sep. 2005.



Jyrki Kivinen received the M.Sc. and Ph.D. degrees in computer science from the University of Helsinki, Finland, in 1989 and 1992, respectively.

He has held various teaching and research appointments at the University of Helsinki and visited the University of California at Santa Cruz and Australian National University as a Postdoctoral Fellow. Since 2003, he has been a Professor at the University of Helsinki. His scientific interests include machine learning and algorithms theory.



Manfred K. Warmuth received the undergraduate degree in computer science from Friedrich Alexander Universität, Germany, in 1978. He received the M.S. and Ph.D. degrees in computer science from the University of Colorado at Boulder in 1980 and 1981, respectively.

Since then he has spent most of his time teaching at the University of California in Santa Cruz. His current research interests are machine learning, online learning, statistical decision theory, and game theory.

Dr. Warmuth received a Fulbright fellowship.



Babak Hassibi was born in Tehran, Iran, in 1967. He received the B.S. degree from the University of Tehran in 1989 and the M.S. and Ph.D. degrees from Stanford University, Stanford, CA, in 1993 and 1996, respectively, all in electrical engineering.

From October 1996 to October 1998, he was a Research Associate with the Information Systems Laboratory, Stanford University. From November 1998 to December 2000, he was a Member of Technical Staff with the Mathematical Sciences Research Center, Bell Laboratories, Murray Hill, NJ. Since January 2001, he has been with the Department of Electrical Engineering, California Institute of Technology, Pasadena, where he is currently an Associate Professor. He has also held short-term appointments at Ricoh California Research Center, the Indian Institute of Science, and Linköping University, Sweden. His research interests include wireless communications, robust estimation and control, adaptive signal processing, and linear algebra. He is the coauthor of *Indefinite Quadratic Estimation and Control: A Unified Approach to H^2 and H^∞ Theories* (New York: SIAM, 1999) and *Linear Estimation* (Englewood Cliffs, NJ: Prentice Hall, 2000).

Prof. Hassibi received an Alborz Foundation Fellowship, the 1999 O. Hugo Schuck Best Paper Award from the American Automatic Control Council, the 2002 National Science Foundation Career Award, the 2002 Okawa Foundation Research Grant for Information and Telecommunications, the 2003 David and Lucille Packard Fellowship for Science and Engineering, and the 2003 Presidential Early Career Award for Scientists and Engineers. He was a Guest Editor for the IEEE TRANSACTIONS ON INFORMATION THEORY special issue on "Space-Time Transmission, Reception, Coding and Signal Processing" and is currently an Associate Editor for Communications of the IEEE TRANSACTIONS ON INFORMATION THEORY.