

---

# The $p$ -norm generalization of the LMS algorithm for adaptive filtering

Jyrki Kivinen  
University of Helsinki

Manfred Warmuth  
University of California, Santa Cruz

Babak Hassibi  
California Institute of Technology

---

## Least Mean Squares (LMS) update

Pick learning rate  $\eta > 0$ . Initialize  $w_0 = \mathbf{0} \in \mathbf{R}^n$

At time  $t$ , for  $t = 1, \dots, T$ , the algorithm

- observes input  $x_t \in \mathbf{R}^n$
- makes prediction  $w_{t-1} \cdot x_t \in \mathbf{R}$
- observes feedback  $y_t \in \mathbf{R}$ , and
- updates its hypothesis as

$$w_t = w_{t-1} - \eta(w_{t-1} \cdot x_t - y_t)x_t$$

## Main Results

- Techniques from machine learning lead to generalizations of LMS
- $H^\infty$ -optimal filtering in signal processing is similar to relative on-line loss bounds in machine learning

## Motivation

- Non-Gaussian modeling
- Get away from rotation invariant algorithms
- Develop algorithms that work well when instances orthogonal and target weight vectors “sparse”

## Expected Bounds for LMS

- Assume  $y_t = \mathbf{u} \cdot \mathbf{x}_t + \nu_t$ , where  $\nu_t$  iid with  $E[\nu_t^2] = \varepsilon$ . Then

$$E \left[ \frac{1}{T} \sum_{t=1}^T (\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 \right] \leq \varepsilon + \frac{1}{T} X_2^2 \|\mathbf{u}\|_2^2$$

Better algorithms exist for probabilistic setting  
However our goal is to weaken the assumptions

## $H^\infty$ bound for LMS [HSK96]

Assume  $\|\mathbf{x}_t\|_2 \leq X_2$  for all  $t$ . Choose  $\eta = 1/X_2^2$   
For any  $\mathbf{u} \in \mathbf{R}^n$

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

- If some  $\mathbf{u}$  with small norm is good predictor then LMS must approximate predictions of  $\mathbf{u}$
- Bound holds for **any**  $\mathbf{u}$  and  $(\mathbf{x}_t, y_t)$
- No probabilistic assumptions
- LMS is  $H^\infty$ -optimal:  
No algorithm can achieve ratio  $< 1$   
 $\forall \mathbf{u}$  and  $(\mathbf{x}_t, y_t)$

## Two related problems

**A priori filtering:** Control Theory

Try to match  $\mathbf{u} \cdot \mathbf{x}_t$

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

**Prediction:** On-line Learning

Try to match  $y_t$

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

## Comparison of known LMS-related bounds

- For  $\eta = \alpha/X_2^2$ ,

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

- For  $\eta = \alpha/X_2^2$  ( $0 < \alpha < 1$ )

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

$$\stackrel{\text{tuned } \alpha}{\leq} Loss_u + 2\sqrt{Loss_u} X_2 \|\mathbf{u}\|_2 + X_2^2 \|\mathbf{u}\|_2^2$$

[CBLW96]

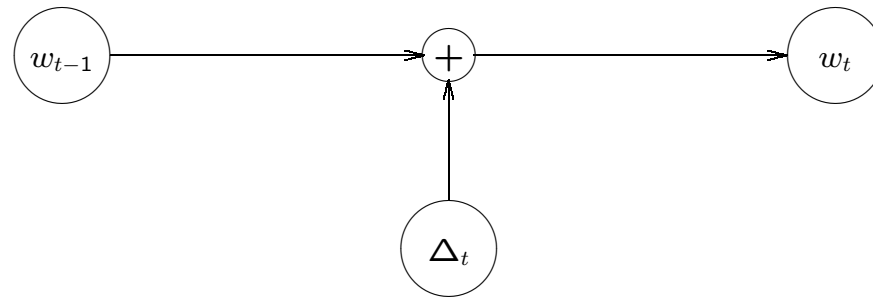


## Generalizing the LMS bound

- Replace  $\|\mathbf{x}\|_2\|\mathbf{u}\|_2$  by  $\|\mathbf{x}\|_p\|\mathbf{u}\|_q$  where  $1/p + 1/q = 1$  and  $\|\mathbf{x}\|_p = (\sum_i |x_i|^p)^{1/p}$
- Instead of comparing predictions to  $\mathbf{u} \cdot \mathbf{x}_t$  for a fixed target  $\mathbf{u}$  compare to  $\mathbf{u}_t \cdot \mathbf{x}_t$  where  $\mathbf{u}_t$  may change
- Replace  $(\cdot)^2$  by more general loss

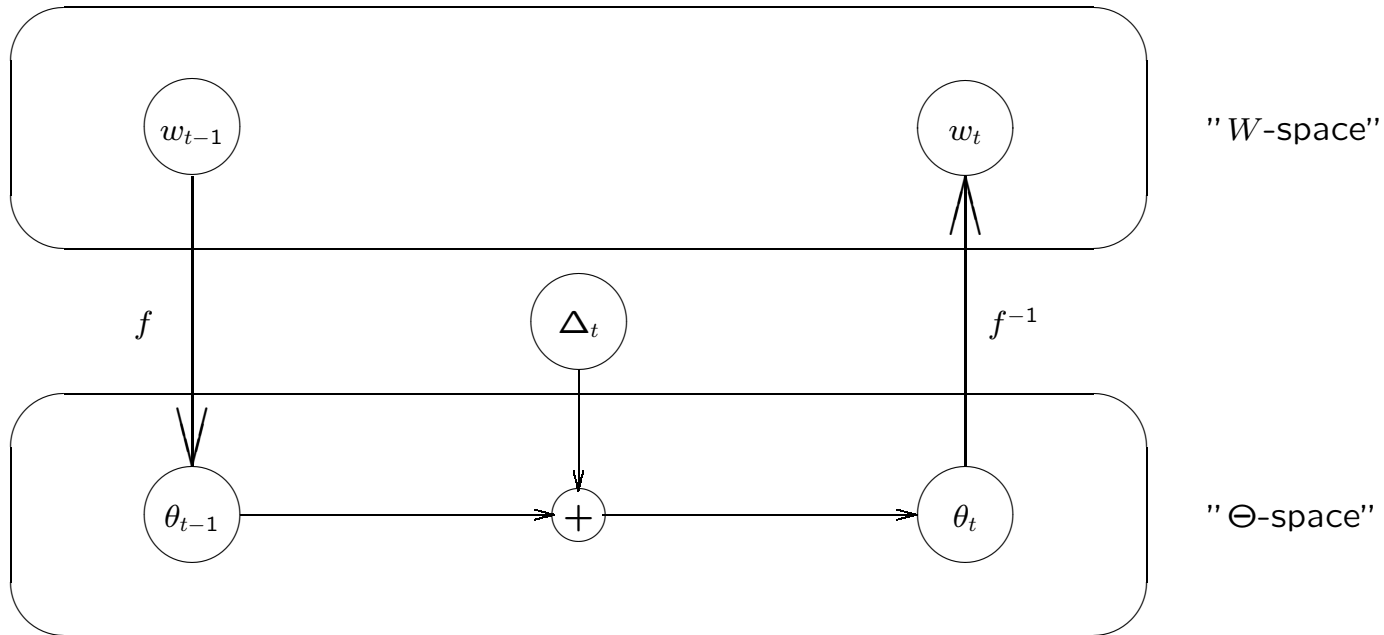
## Basic LMS

$$\Delta_t = -\eta(w_{t-1} \cdot x_t - y_t)x_t$$



## $p$ -norm LMS

Write  $\theta_t = f(w_t)$



**$p$ -norm LMS** based on [GLS01]

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

where  $\mathbf{f}: \mathbf{R}^n \rightarrow \mathbf{R}^n$  given by

$$f_i(\mathbf{w}) = \frac{\text{sign}(w_i)|w_i|^{q-1}}{\|\mathbf{w}\|_q^{q-2}} \quad \text{and} \quad f_i^{-1}(\boldsymbol{\theta}) = \frac{\text{sign}(\theta_i)|\theta_i|^{p-1}}{\|\boldsymbol{\theta}\|_p^{p-2}}$$

When  $p = q = 2$ , then  $\mathbf{f}(\mathbf{w}) = \mathbf{w}$ : LMS

For large  $p$ ,  $\mathbf{f}^{-1}$  emphasizes differences in components

## A priori filtering bound

**Theorem** Assume  $\|\mathbf{x}_t\|_p \leq X_p$  for all  $t$ , and let  $\eta = 1/((p-1)X_p^2)$ . Then for any  $\mathbf{u}$  the  $p$ -norm algorithm satisfies

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

- $1/p + 1/q = 1$  and  $2 \leq p < \infty$ ,  $1 < q \leq 2$
- How do we get the dual norm pair  $(\infty, 1)$  (where  $\|\mathbf{x}\|_\infty = \max_i |x_i|$ )?

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

## Comparison with basic LMS

New bounds incomparable with old ones because for  $p > 2$  and  $q < 2$

$$\|\mathbf{x}\|_p < \|\mathbf{x}\|_2 \quad \text{and} \quad \|\mathbf{u}\|_q > \|\mathbf{u}\|_2$$

Compare  $p = 2$  and  $p = O(\log n)$  in two extreme cases:

**Sparse target, dense instances:** Let  $\mathbf{u} = (1, 0, \dots, 0)$  and  $\mathbf{x} = (1, \dots, 1)$ .

- $\|\mathbf{x}\|_2^2 \|\mathbf{u}\|_2^2 = n^2$
- $(\log n) \|\mathbf{x}\|_\infty^2 \|\mathbf{u}\|_1^2 = \log n$
- Thus large  $p$  better

**Dense target, sparse instances:** Let  $\mathbf{u} = (1, 1, \dots, 1)$  and  $\mathbf{x} = (1, 0, \dots, 0)$ .

- $\|\mathbf{x}\|_2^2 \|\mathbf{u}\|_2^2 = n^2$
- $(\log n) \|\mathbf{x}\|_\infty^2 \|\mathbf{u}\|_1^2 = n^2 \log n$
- Thus  $p = 2$  better

## The $p$ -norm LMS can behave like EG

### Hadamard Matrix:

	→	+1	+1	+1	+1
<i>instances</i>	→	+1	-1	+1	-1
	→	+1	+1	-1	-1
	→	+1	-1	-1	+1
		↑	↑	↑	↑
			<i>targets</i>		

- Instances are orthogonal
- Target weight vectors are units
- LMS: error  $\geq 1 - \frac{k}{n}$
- $p$ -norm LMS with  $p = O(\log n)$ : error  $\geq \frac{\ln n}{k}$

## Time-varying target (following [HW01])

Up to now, model has been  $y_t = \mathbf{u} \cdot \mathbf{x}_t + \text{noise}$  where target  $\mathbf{u}$  is fixed

Generalize this to  $y_t = \mathbf{u}_t \cdot \mathbf{x}_t + \text{noise}$  where target  $\mathbf{u}_t$  may vary over time

**Example 1: target makes one jump** Choose  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$  and take

$$\mathbf{u}_t = \begin{cases} \mathbf{a} & \text{for } 1 \leq t \leq T/2 \\ \mathbf{b} & \text{for } T/2 < t \leq T \end{cases}$$

**Example 2: target moves steadily** Choose  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$  and take

$$\mathbf{u}_t = \frac{T-t}{T-1} \mathbf{a} + \frac{t-1}{T-1} \mathbf{b}$$



## Algorithms for time-varying target

Old update:

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

Bounding update:

$$\mathbf{w}_t = \begin{cases} \mathbf{w}'_t & \text{if } \|\mathbf{w}'_t\|_q \leq U_q \\ U_q \frac{\mathbf{w}'_t}{\|\mathbf{w}'_t\|_q} & \text{otherwise} \end{cases}$$

where  $U_q > 0$  is a norm bound

We rescale the weight vector whenever  $q$ -norm larger than  $U_q$

## Bound for time-varying target

**Theorem** Assume  $\|\mathbf{x}_t\|_p \leq X_p$  for all  $t$ , and let  $\eta = 1/((p-1)X_p^2)$ . Then if  $\|\mathbf{u}_t\|_q \leq U_q$  for all  $t$ , the bounded  $p$ -norm LMS satisfies

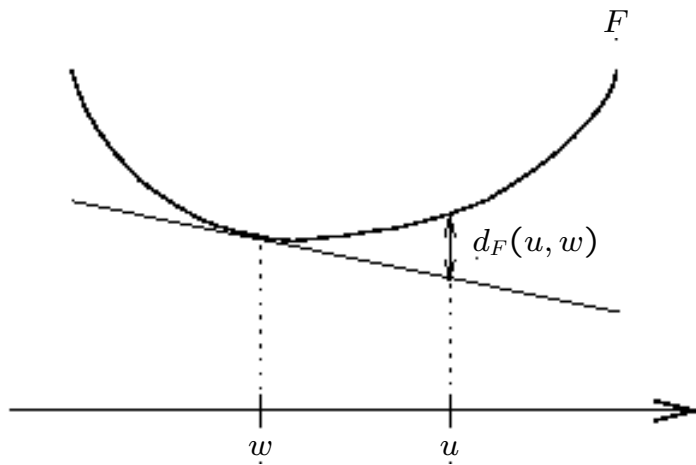
$$\begin{aligned} \sum_{t=1}^T (\mathbf{u}_t \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2 &\leq \sum_{t=1}^T (y_t - \mathbf{u}_t \cdot \mathbf{x}_t)^2 + (p-1)X_p^2 U_q^2 \\ &\quad + 2(p-1)X_p^2 U_q \sum_{t=1}^{T-1} \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q \end{aligned}$$

- Only total distance  $\sum_t \|\mathbf{u}_{t+1} - \mathbf{u}_t\|_q$  traveled by the target matters
- Cost  $2(p-1)X_p^2 U_q$  per unit target movement
- For fixed target  $\mathbf{u}_{t+1} = \mathbf{u}_t$ , we recover previous bound
- However  $U_q$  needs to be known in advance

## Bregman divergences

- Key tool in analyzing and understanding the algorithms
- Fix strictly convex differentiable  $F: \mathbf{R}^n \rightarrow \mathbf{R}$ .  
Denote the gradient by  $\mathbf{f} = \nabla F$ .
- Now the Bregman divergence  $d_F: \mathbf{R}^n \times \mathbf{R}^n \rightarrow \mathbf{R}$  is

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



$d_F(\mathbf{u}, \mathbf{w})$  is the error of first-order Taylor approximation of  $F(\mathbf{u})$  around  $\mathbf{w}$

## Basic properties of Bregman divergences

- $d_F(\mathbf{u}, \mathbf{w}) \geq 0$ ,  $d_F(\mathbf{u}, \mathbf{w}) = 0$  iff  $\mathbf{u} = \mathbf{w}$
- **not** symmetrical (in general)
- does **not** satisfy triangle inequality
- $d_F(\mathbf{u}, \mathbf{w})$  convex in  $\mathbf{u}$ , not necessarily in  $\mathbf{w}$

Connection to **exponential families** (roughly):

- $F$  is cumulant function,  $f$  is link function
- $\mathbf{w}$  is expectation parameter,  $f(\mathbf{w})$  canonical parameter
- $d_F(\mathbf{u}, \mathbf{w})$  is the KL divergence between distributions parameterized by  $\mathbf{u}$  and  $\mathbf{w}$

## Example: $p$ -norm divergence [GLS01]

- 

$$F(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_q^2$$

- Then the gradient  $\mathbf{f} = \nabla F$  satisfies

$$f_i(\mathbf{w}) = \frac{\text{sign}(w_i) |w_i|^{q-1}}{\|\mathbf{w}\|_q^{q-2}} \quad \text{and} \quad f_i^{-1}(\boldsymbol{\theta}) = \frac{\text{sign}(\theta_i) |\theta_i|^{p-1}}{\|\boldsymbol{\theta}\|_p^{p-2}}$$

- The divergence is

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

- Special case  $p = q = 2$  gives  $d_F(\mathbf{u}, \mathbf{w}) = \frac{1}{2} \|\mathbf{u} - \mathbf{w}\|_2^2$

## Deriving the updates

- Define a regularized instantaneous loss

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

- Basic aim is to have

$$\mathbf{w}_t = \underset{\mathbf{w}}{\operatorname{argmin}} C_t(\mathbf{w})$$

- Minimize by setting  $\nabla C_t(\mathbf{w}_t) = \mathbf{0}$ , obtaining the **implicit update**

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

- Approximate  $\mathbf{w}_t \cdot \mathbf{x}_t \approx \mathbf{w}_{t-1} \cdot \mathbf{x}_t$  to obtain the **update**

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

## Analyzing the update

- Measure of progress

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

- Massage the term  $(y_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t) \mathbf{x}_t \cdot (\mathbf{u} - \mathbf{w}_{t-1})$  until  $(\mathbf{u} \cdot \mathbf{x}_t - \mathbf{w}_{t-1} \cdot \mathbf{x}_t)^2$  and  $(y_t - \mathbf{u} \cdot \mathbf{x}_t)^2$  appear; throw rest away
- Estimate  $d_F(\mathbf{w}_{t-1}, \mathbf{w}_t)$  in terms of  $\|\mathbf{x}_t\|_p$
- Sum over  $t = 1, \dots, T$

## Conclusion

- LMS and normalized LMS can be derived from an optimization problem involving a certain Bregman divergence
- Different Bregman divergences lead to different algorithms, with loss bounds in terms of different norms
- Bounds can be generalized for time-varying targets (and generalized linear models, not presented in the talk); proofs easy
- Algorithms for  $p = 2$  can be kernelized, for  $p > 2$  probably not

**Bottom line:** Machinery from on-line machine learning carries over to  $H^\infty$ -optimal filtering



## Where are we headed?

- Develop  $p$ -norm Kalman filter
- Prove relative loss bounds