

Machine Learning Algorithms

Perceptron Learning Algorithm (PLA)

Markus Rupp

12.9.2020



Recall: Binary Classification

- Binary classification: y from $\{-1, 1\}$

$$\underline{\hat{x}}_i^H \underline{\hat{c}} > 0; \quad \text{if } y_i = +1$$

$$\underline{\hat{x}}_i^H \underline{\hat{c}} < 0; \quad \text{if } y_i = -1$$

- We like to find an (iterative) algorithm that learns from data pairs and finds the hyperplane separating both categories



Sigmoid or tanh?

- For this we employ a nonlinear function $f(x)$
 - That maps the linear input (vector product) to one of the classes
 - And is sufficiently smooth, i.e., differentiable

$$\sigma(\underline{\hat{x}}_i^H \underline{c}) = y_i + v_i$$

$$\sigma(x) = \frac{1}{1 + e^{-\beta x}} \in [0, 1]$$

$$\begin{aligned} 2\sigma(x) - 1 &= \frac{2}{1 + e^{-\beta x}} - 1 = \frac{1 - e^{-\beta x}}{1 + e^{-\beta x}} \\ &= \tanh\left(\frac{\beta x}{2}\right) \in [-1, 1] \end{aligned}$$

$$\frac{\partial \sigma(x)}{\partial x} = \beta \frac{e^{-\beta x}}{(1 + e^{-\beta x})^2} = \beta \sigma(x)(1 - \sigma(x))$$

$$\frac{\partial \tanh\left(\frac{\beta x}{2}\right)}{\partial x} = 2\beta \sigma(x)(1 - \sigma(x)) = \frac{\beta}{2} \left(1 - \tanh^2\left(\frac{\beta x}{2}\right)\right)$$



PLA

- In principle both are valid options
- Tanh is usually preferred due to symmetry
- For a function $f(x)$ (sigmoid or tanh), we obtain the
- Perceptron Learning Algorithm (PLA)

$$\underline{\hat{W}}_k = \underline{\hat{W}}_{k-1} + \mu \underline{\mathbf{x}}_k^* \left(y_k - f \left(\underline{\mathbf{x}}_k^T \underline{\hat{W}}_{k-1} \right) \right); k = 1, 2, \dots$$



PLA Derivation

- Let's derive this from first principles. We start with the cost function

$$c(\underline{w}) = \left(y_k - f\left(\underline{x}_k^T \underline{w}\right) \right)^2$$

- Its derivative is then given by

$$\frac{\partial}{\partial \underline{w}} c(\underline{w}) = -2 \left(y_k - f\left(\underline{x}_k^T \underline{w}\right) \right) f'\left(\underline{x}_k^T \underline{w}\right) \underline{x}_k = \underbrace{-2e(\underline{w})\underline{x}_k}_{\text{classical LS term}} f'\left(\underline{x}_k^T \underline{w}\right)$$



PLA Derivation

- Using such gradient we find

$$\underline{\hat{\mathbf{w}}}_k = \underline{\hat{\mathbf{w}}}_{k-1} + \tilde{\mu} f' \left(\underline{\mathbf{x}}_k^T \underline{\hat{\mathbf{w}}}_{k-1} \right) \underline{\mathbf{x}}_k^* \left(y_k - f \left(\underline{\mathbf{x}}_k^T \underline{\hat{\mathbf{w}}}_{k-1} \right) \right); k = 1, 2, \dots$$

- Since the sigmoid function ends up in a limited (positive) range of potential values, we consume the derivative into the new step-size μ and simplify the algorithmic update to obtain:

$$\underline{\hat{\mathbf{w}}}_k = \underline{\hat{\mathbf{w}}}_{k-1} + \mu \underline{\mathbf{x}}_k^* \left(y_k - f \left(\underline{\mathbf{x}}_k^T \underline{\hat{\mathbf{w}}}_{k-1} \right) \right); k = 1, 2, \dots$$



Analysis

- Unfortunately, this small modification of the LMS algorithm leads to an unexpected hurdle in analysing its behaviour.
- Our classical form of analysis cannot be applied any longer and a new method needs to be applied.
- 1) show new method (based on energy relations) on LMS algorithm
- 2) apply it for PLA
- 3) apply it to recurrent neural nets (RNN)



Energy based relations

- As the entire analysis is based on the Cauchy Schwarz inequality, let's recall it....



Cauchy-Schwarz' Inequality

- **Theorem:** In an inner vector product space S in C^n with induced norm $||\cdot||$ the following holds:

$$|\langle \underline{x}, \underline{y} \rangle| \leq \|\underline{x}\|_2 \|\underline{y}\|_2$$

- The inequality holds with equal sign, if and only if, $\underline{x} = \alpha \underline{y}$.
- **Proof:** We start with the simple fact that

$$\|\underline{x} - \alpha \underline{y}\|_2 \geq 0$$



Cauchy-Schwarz' Inequality

- Furthermore, we have

$$\begin{aligned} 0 &\leq \|\underline{x} - \alpha \underline{y}\|_2^2 = \|\underline{x}\|_2^2 - 2 \operatorname{Re} \langle \underline{x}, \alpha \underline{y} \rangle + |\alpha|^2 \|\underline{y}\|_2^2 \\ &= \|\underline{x}\|_2^2 - \frac{|\langle \underline{x}, \underline{y} \rangle|^2}{\|\underline{y}\|_2^2} + \|\underline{y}\|_2^2 \left[\left(\alpha - \frac{\langle \underline{x}, \underline{y} \rangle}{\|\underline{y}\|_2^2} \right) \left(\alpha^* - \frac{\langle \underline{x}, \underline{y} \rangle^*}{\|\underline{y}\|_2^2} \right) \right] \end{aligned}$$

- The minimum is obtained for:

$$\alpha = \frac{\langle \underline{x}, \underline{y} \rangle}{\|\underline{y}\|_2^2}$$

- Thus:

$$0 \leq \min_{\alpha} \|\underline{x} - \alpha \underline{y}\|_2^2 = \|\underline{x}\|_2^2 - \frac{|\langle \underline{x}, \underline{y} \rangle|^2}{\|\underline{y}\|_2^2} \rightarrow |\langle \underline{x}, \underline{y} \rangle| \leq \|\underline{x}\|_2 \|\underline{y}\|_2$$



Cauchy-Schwarz' Inequality

- We are not done yet!
- For $\underline{x} = \alpha \underline{y}$ we obtain $||\underline{x} - \alpha \underline{y}|| = 0$.
- Due to its norm property $||\underline{x}|| \geq 0$, this property $||\underline{x} - \alpha \underline{y}|| = 0$ can only be achieved for the argument to be zero.
- Thus, we obtain equality if and only if $\underline{x} = \alpha \underline{y}$.



Local Passivity Properties

- Now consider the term $\mu_k \|\underline{x}_k\|^2 \leq 1$
- Take an arbitrary estimate \underline{q} .
- With Cauchy-Schwarz' inequality:

$$\frac{\left| \underline{x}_k^T \underline{w} - \underline{x}_k^T \underline{q} \right|^2}{\mu_k^{-1} \left\| \underline{w} - \underline{q} \right\|_2^2} \leq 1$$

- Alternatively:

$$\left| \underline{x}_k^T \underline{w} - \underline{x}_k^T \underline{q} \right|^2 - \mu_k^{-1} \left\| \underline{w} - \underline{q} \right\|_2^2 \leq 0$$



Local Passivity Properties

- Such relation is still valid, once we add something positive in the denominator

$$\frac{\left| \underline{x}_k^T \underline{w} - \underline{x}_k^T \underline{q} \right|^2}{\mu_k^{-1} \left\| \underline{w} - \underline{q} \right\|_2^2 + \left| v_k \right|^2} \leq 1$$

- Even for particular estimates, this is true:

$$\frac{\left| \underline{x}_k^T \underline{w} - \underline{x}_k^T \hat{\underline{w}}_{k-1} \right|^2}{\mu_k^{-1} \left\| \underline{w} - \hat{\underline{w}}_{k-1} \right\|_2^2 + \left| v_k \right|^2} \leq 1$$



Local Passivity Relations

- We are facing now the question, whether some algorithms with their specific estimates, like LMS or RLS, (or others) improve such relation even further by lowering the limit below one.



Local Passivity Properties

Theorem 5.1 (Local Passivity) For the adaptive gradient-type algorithm (LMS with variable step-size) for every time-instant the following properties are true:

$$a) \quad \frac{\mu_k^{-1} \|\underline{w} - \hat{\underline{w}}_k\|_2^2 + |e_{a,k}|^2}{\mu_k^{-1} \|\underline{w} - \hat{\underline{w}}_{k-1}\|_2^2 + |v_k|^2} \leq 1$$

$$b) \quad \frac{|e_{p,k}|^2 + |e_{a,k}|^2}{\mu_k^{-1} \|\underline{w} - \hat{\underline{w}}_{k-1}\|_2^2 + |v_k|^2} \leq 1$$

$$c) \quad \frac{\gamma_k \|\underline{w} - \hat{\underline{w}}_k\|_2^2 + |e_{p,k}|^2}{\gamma_k \|\underline{w} - \hat{\underline{w}}_{k-1}\|_2^2 + |v_k|^2} \leq 1$$



Local Passivity Properties

$$d) \quad \frac{|e_{a,k}|^2 + |e_{a,k+1}|^2}{\mu_k^{-1} \|\underline{w} - \hat{\underline{w}}_{k-1}\|_2^2 + |\underline{v}_k|^2} \leq 1$$

**The first three relations hold for $\mu_k \|\underline{x}_k\|_2^2 \leq 1$ while
for the last relation $\mu_k \leq \min \left\{ 1 / \|\underline{x}_k\|_2^2, 1 / \|\underline{x}_{k+1}\|_2^2 \right\}$**

Also, $\gamma_k = \mu_k^{-1} - \|\underline{x}_k\|_2^2$



Local Passivity Relations

- **Proof:** Starting with the first relation, the update equation of the gradient type algorithm reads:

$$\underline{\tilde{w}}_k = \underline{\tilde{w}}_{k-1} - \mu_k \underline{x}_k^* [e_{a,k} + v_k]$$

Here, the disturbed error signal was split in the undisturbed part and the additive noise:

$$\tilde{e}_{a,k} = e_{a,k} + v_k$$



Local Passivity Relations

- Computing the squared L_2 norm of both sides of the equation, we obtain:

$$\begin{aligned}\|\underline{\tilde{w}}_k\|_2^2 &= \|\underline{\tilde{w}}_{k-1}\|_2^2 + \mu_k^2 \|\underline{x}_k\|_2^2 |e_{a,k} + v_k|^2 \\ &\quad - \mu_k [e_{a,k} + v_k] e_{a,k}^* - \mu_k [e_{a,k}^* + v_k^*] e_{a,k} \\ |e_{a,k} + v_k|^2 &= |e_{a,k}|^2 + |v_k|^2 + v_k e_{a,k}^* + v_k^* e_{a,k}\end{aligned}$$

- And thus:

$$\begin{aligned}\|\underline{\tilde{w}}_k\|_2^2 &= \|\underline{\tilde{w}}_{k-1}\|_2^2 + \mu_k^2 \|\underline{x}_k\|_2^2 |e_{a,k} + v_k|^2 \\ &\quad - 2\mu_k |e_{a,k}|^2 - \mu_k [|e_{a,k} + v_k|^2 - |e_{a,k}|^2 - |v_k|^2]\end{aligned}$$



Local Passivity Relations

- By reordering we finally obtain:

$$\begin{aligned} & \|\underline{\tilde{w}}_k\|_2^2 - \|\underline{\tilde{w}}_{k-1}\|_2^2 + \mu_k |e_{a,k}|^2 - \mu_k |v_k|^2 \\ &= \mu_k |e_{a,k} + v_k|^2 \left[\mu_k \|\underline{x}_k\|_2^2 - 1 \right] \end{aligned}$$

- As long as $\mu_k \|\underline{x}_k\|^2 \leq 1$ the right hand side remains negative, and thus

$$\|\underline{\tilde{w}}_k\|_2^2 - \|\underline{\tilde{w}}_{k-1}\|_2^2 + \mu_k |e_{a,k}|^2 - \mu_k |v_k|^2 \leq 0$$

- Division gives us the desired result.



Robustness Analysis for LMS

- Need global relation from $k=1,2,\dots,N$

$$-\|\underline{\tilde{w}}_k\|_2^2 + \|\underline{\tilde{w}}_{k-1}\|_2^2 + \mu_k |v_k|^2 \geq \mu_k |e_{a,k}|^2$$

- Sum up and obtain:

$$\frac{\|\underline{\tilde{w}}_N\|_2^2 + \sum_{k=1}^N \mu_k |e_{a,k}|^2}{\|\underline{\tilde{w}}_0\|_2^2 + \sum_{k=1}^N \mu_k |v_k|^2} \leq 1$$



Robustness Analysis for LMS

- Find passivity relation as well:

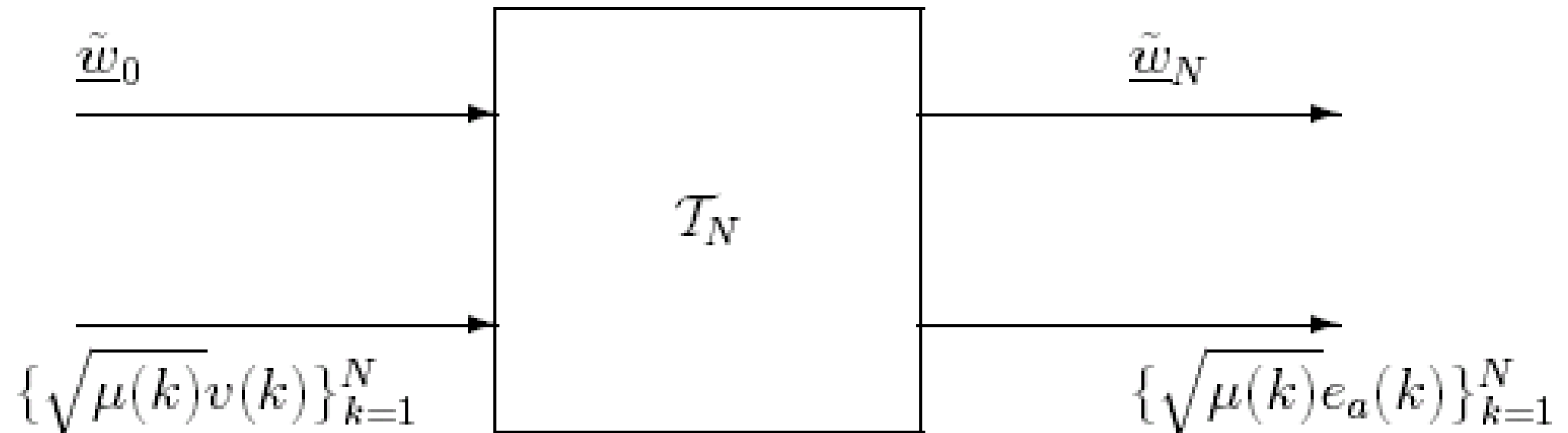
$$\begin{bmatrix} \sqrt{\mu_1} e_{a,1} \\ \vdots \\ \sqrt{\mu_N} e_{a,N} \\ \underline{\tilde{w}}_N \end{bmatrix} = \underbrace{\begin{bmatrix} x & & & \\ x & x & & \\ & x & x & \\ x & x & x & x \end{bmatrix}}_{\mathbf{T}_N} \begin{bmatrix} \underline{\tilde{w}}_0 \\ \sqrt{\mu_1} v_1 \\ \vdots \\ \sqrt{\mu_1} v_1 \end{bmatrix}$$

- Induced matrix norm is bounded:

$$\|\mathbf{T}_N\|_{2,ind} \leq 1$$



Robustness Analysis for LMS



Minimax Optimality

- Is there a sequence so that the maximum one is really obtained?

$$\frac{\|\underline{\tilde{w}}_N\|_2^2 + \sum_{k=1}^N \mu_k |e_{a,k}|^2}{\|\underline{\tilde{w}}_0\|_2^2 + \sum_{k=1}^N \mu_k |v_k|^2} \stackrel{?}{=} 1$$

Select special noise sequence

$$v_k = -e_{a,k}$$



Minimax Optimality

- However, in such case no update will occur, and the final estimate is identical to the first one, thus

$$\max_{\hat{\underline{w}}_0 \neq \underline{w}, v(\cdot)} \frac{\|\underline{w} - \hat{\underline{w}}_N\|_2^2 + \sum_{k=1}^N \mu_k |e_{a,k}|^2}{\|\underline{w} - \hat{\underline{w}}_0\|_2^2 + \sum_{k=1}^N \mu_k |v_k|^2} = 1$$



Minimax Optimality

- On the other hand, for each algorithm we have for $v_k = -e_{a,k}$ and

$$: \quad \sum_{k=1}^N \mu_k |v_k|^2 = \sum_{k=1}^N \mu_k |e_{a,k}|^2 \leq \|\underline{w} - \hat{\underline{w}}_N\|_2^2 + \sum_{k=1}^N \mu_k |e_{a,k}|^2$$

- or equivalently (for $\hat{\underline{w}}_0 = \underline{w}$): $\underline{w} = \hat{\underline{w}}_0$

$$\max_{v(.)} \frac{\|\underline{w} - \hat{\underline{w}}_N\|_2^2 + \sum_{k=1}^N \mu_k |e_{a,k}|^2}{\|\underline{w} - \hat{\underline{w}}_0\|_2^2 + \sum_{k=1}^N \mu_k |v_k|^2} \geq 1$$



Minimax Optimality

Theorem 5.2: (Minimax Property of the Gradient Type Algorithm) The gradient type algorithm solves the following Minimax problem if

$$\mu_k \|\underline{x}_k\|^2 \leq 1$$

$$\min_{\text{All algorithms}} \max_{\tilde{w}_0, v(\cdot)} \frac{\|\underline{w} - \hat{w}_N\|_2^2 + \sum_{k=1}^N \mu_k |e_{a,k}|^2}{\|\underline{w} - \hat{w}_0\|_2^2 + \sum_{k=1}^N \mu_k |v_k|^2}$$

The optimal value is one.



Convergence

- **Theorem 5.3 (Convergence Conditions)**

$$\text{If } \mu_k \|\underline{x}_k\|^2 \leq 1, \|\underline{w}_0\| < \infty, \text{ and } \sum_{k=1}^{\infty} \mu_k |v_k|^2 < \infty$$

then it follows that: $\sqrt{\mu_k} e_{a,k} \rightarrow 0$.

If, furthermore, $\sqrt{\mu_k} \underline{x}_k$ is persistent exciting, then we also have:

$$\underline{\hat{w}}_k \rightarrow \underline{w}$$



Convergence

- **Proof:** Since we have

$$\sum_{k=1}^N \mu_k |e_{a,k}|^2 \leq \|\underline{w} - \hat{w}_0\|_2^2 + \sum_{k=1}^N \mu_k |v_k|^2$$

- If initial error as well as noise energy is limited, then the a-priori error energy must be limited and thus must decay after some time to zero.
- Cauchy-Series!



Convergence

- Consider the update equation at time instant k:

$$\underline{\tilde{w}}_k = \underline{\tilde{w}}_{k-1} - \mu_k \underline{x}_k^* [e_{a,k} + v_k]$$

- Summing up over p time instants:

$$\underline{\tilde{w}}_{k+p-1} = \underline{\tilde{w}}_k - \sum_{l=1}^{p-1} \mu_{k+l} \tilde{e}_{a,k+l} \underline{x}_{k+l}^*$$

- Consider a series of $P > M$ vectors \underline{x}_{k+p} , $p=1, \dots, P$, all such vectors can be tested at one value of the parameter error vector:

- $$\underline{x}_{k+p}^T \underline{\tilde{w}}_k = \underbrace{\underline{x}_{k+p}^T \underline{\tilde{w}}_{k+p-1}}_{e_{a,k+p}} - \sum_{l=1}^{p-1} \mu_{k+l} \tilde{e}_{a,k+l} \underline{x}_{k+p}^T \underline{x}_{k+l}^*$$



Convergence

- Since all $\underline{e}_{a,k} \rightarrow 0$ and since also all $\underline{v}_k \rightarrow 0$, we also have that
$$\tilde{\underline{e}}_{a,k} \rightarrow 0$$
- Since the right hand side goes to zero, the left hand side must tend to zero as well. Since this is true for all vectors, we can write:

$$\begin{bmatrix} \underline{x}_{k+1}^T \\ \underline{x}_{k+2}^T \\ \vdots \\ \underline{x}_{k+P}^T \end{bmatrix} \tilde{\underline{w}}_k \rightarrow 0$$



Convergence

- From here we cannot conclude yet that $\underline{\tilde{w}}_k \rightarrow \underline{0}$ since the vector could exist in the null space of \underline{x}_k .
- Invoking the persistent excitation condition as well, the matrix must have the full rank M, or equivalently:

$$\begin{bmatrix} \underline{x}_{k+1}^* & \underline{x}_{k+2}^* & \dots & \underline{x}_{k+P}^* \end{bmatrix} \begin{bmatrix} \underline{x}_{k+1}^T \\ \underline{x}_{k+2}^T \\ \vdots \\ \underline{x}_{k+P}^T \end{bmatrix} \underline{\tilde{w}}_k \rightarrow 0$$



Convergence

- Since the persistent excitation condition requires that

$$0 < \alpha I \leq \begin{bmatrix} \underline{x}_{k+1}^* & \underline{x}_{k+2}^* & \dots & \underline{x}_{k+P}^* \end{bmatrix} \begin{bmatrix} \underline{x}_{k+1}^T \\ \underline{x}_{k+2}^T \\ \vdots \\ \underline{x}_{k+P}^T \end{bmatrix} \leq \beta I$$

- We can only conclude that

$$\underline{\tilde{w}}_k \rightarrow \underline{0}$$



Small Gain Theorem

- Consider the Small Gain Theorem as an extension to what you know about stability of linear systems to those systems that are non-linear.
- Remember in linear system theory, the open loop gain is required to be smaller than one in order to have a closed loop gain bounded (BIBO stability).
- The Small Gain Theorem applies equally for non-linear systems as well as time-variant systems.
- Consider an input and output sequence as entries of a vector:

$$\underline{y}_N = H_N \underline{x}_N$$



Small Gain Theorem

- **Definition 1:** A mapping H is called l -stable if two positive constants β, γ exist, such that for all input sequences the following is true:

$$\|\underline{y}_N\| = \|H_N \underline{x}_N\| \leq \gamma \|\underline{x}_N\| + \beta$$

- **Definition 2:** The smallest positive value γ for which l -stability is guaranteed, is called the gain of H .
- **Remark:** BIBO (Bounded Input Bounded Output) stability is equivalent to l_{oo} - stability.



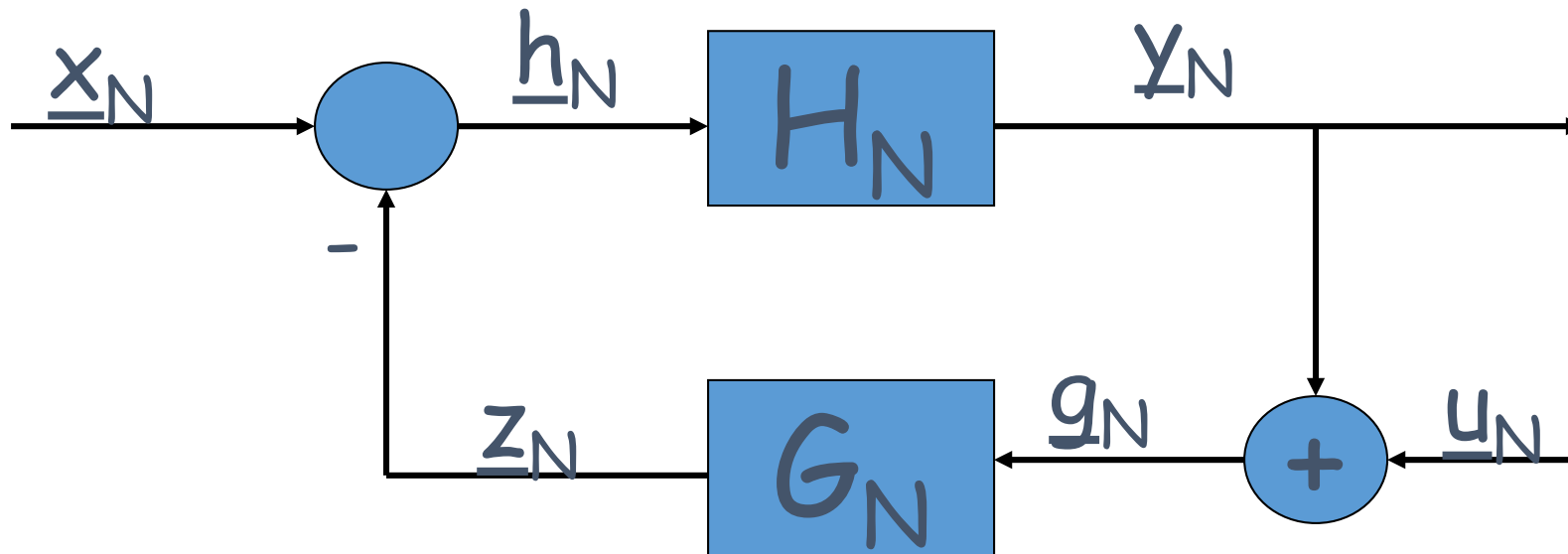
Small Gain Theorem

- Now consider a feedback system comprising of two systems G and H with corresponding gains γ_g and γ_h with the following signals:

$$\underline{y}_N = H_N \underline{h}_N = H_N [\underline{x}_N - \underline{z}_N]$$
$$\underline{z}_N = G_N \underline{g}_N = G_N [\underline{u}_N + \underline{y}_N]$$



Small Gain Theorem



Small Gain Theorem

- **Small Gain Theorem:** If the two gains γ_g and γ_h are such that

$$\gamma_g \gamma_h < 1$$

then the signals \underline{g}_N and \underline{h}_N of the feedback system are bounded by:

$$\|\underline{h}_N\| \leq \frac{1}{1 - \gamma_g \gamma_h} \left[\|\underline{x}_N\| + \gamma_g \|\underline{u}_N\| + \beta_g + \gamma_g \beta_h \right]$$

$$\|\underline{g}_N\| \leq \frac{1}{1 - \gamma_g \gamma_h} \left[\|\underline{u}_N\| + \gamma_h \|\underline{x}_N\| + \beta_h + \gamma_h \beta_g \right]$$



Small Gain Theorem

- **Proof:** Starting with the norm is given by:

$$\underline{h}_N = \underline{x}_N - \underline{z}_N; \quad \underline{g}_N = \underline{u}_N + \underline{y}_N$$

$$\begin{aligned} \|\underline{h}_N\| &\leq \|\underline{x}_N\| + \|G_N \underline{g}_N\| \\ &\leq \|\underline{x}_N\| + \gamma_g \|\underline{g}_N\| + \beta_g \\ &\leq \|\underline{x}_N\| + \gamma_g [\|\underline{u}_N\| + \gamma_h \|\underline{h}_N\| + \beta_h] + \beta_g \\ &= \gamma_g \gamma_h \|\underline{h}_N\| + \|\underline{x}_N\| + \gamma_g \|\underline{u}_N\| + \beta_g + \gamma_g \beta_h \\ &= \frac{1}{1 - \gamma_g \gamma_h} [\|\underline{x}_N\| + \gamma_g \|\underline{u}_N\| + \beta_g + \gamma_g \beta_h] \end{aligned}$$

- The derivation for \underline{g}_N follows a similar path.



Robustness: Feedback Formulation

- For LMS we derived:

$$\begin{aligned} \|\underline{\tilde{w}}_k\|_2^2 - \|\underline{\tilde{w}}_{k-1}\|_2^2 + \mu_k |e_{a,k}|^2 - \mu_k |v_k|^2 \\ = \mu_k |e_{a,k} + v_k|^2 \left[\mu_k \|\underline{x}_k\|_2^2 - 1 \right] \end{aligned}$$

- Leading to three conditions, for

$$\mu_k \|\underline{x}_k\|_2^2 < 1 \rightarrow \textit{passivity}$$

$$\mu_k \|\underline{x}_k\|_2^2 = 1 \rightarrow \textit{allpass}, \mu_k = \bar{\mu}_k = \frac{1}{\|\underline{x}_k\|_2^2}$$

$$\mu_k \|\underline{x}_k\|_2^2 > 1 \rightarrow ?$$



Robustness: Feedback Formulation

- From a (temporary) local passivity relation for LMS

$$\frac{\|\tilde{\underline{w}}_k\|_2^2 + \mu_k |e_{a,k}|^2}{\|\tilde{\underline{w}}_{k-1}\|_2^2 + \mu_k |v_k|^2} \begin{cases} \leq 1 & \text{if } 0 < \mu_k < \bar{\mu}_k \\ = 1 & \text{if } \mu_k = \bar{\mu}_k \\ > 1 & \text{if } \mu_k > \bar{\mu}_k \end{cases}$$

- We can derive a global relation



Feedback

- **Lemma 5.1:** For the gradient type algorithm the following is true at every time instant k (with $\mu_k = 1/\|x_k\|_2^2$):

$$\frac{\|\tilde{w}_k\|_2^2 + \mu_k |e_{a,k}|^2}{\|\tilde{w}_{k-1}\|_2^2 + \mu_k |v_k|^2} \begin{cases} \leq 1 & \text{if } 0 < \mu_k < \bar{\mu}_k \\ = 1 & \text{if } \mu_k = \bar{\mu}_k \\ > 1 & \text{if } \mu_k > \bar{\mu}_k \end{cases}$$

- **Proof:** The first has been proven already. The second can be directly obtained by setting $\mu_k = \bar{\mu}_k$. For the third one, consider again:

$$\|\tilde{w}_k\|_2^2 - \|\tilde{w}_{k-1}\|_2^2 + \mu_k |e_{a,k}|^2 - \mu_k |v_k|^2 = \mu_k |e_{a,k} + v_k|^2 \left[\underbrace{\mu_k \|x_k\|_2^2}_{>1} - 1 \right]$$



Feedback

- Reformulation:

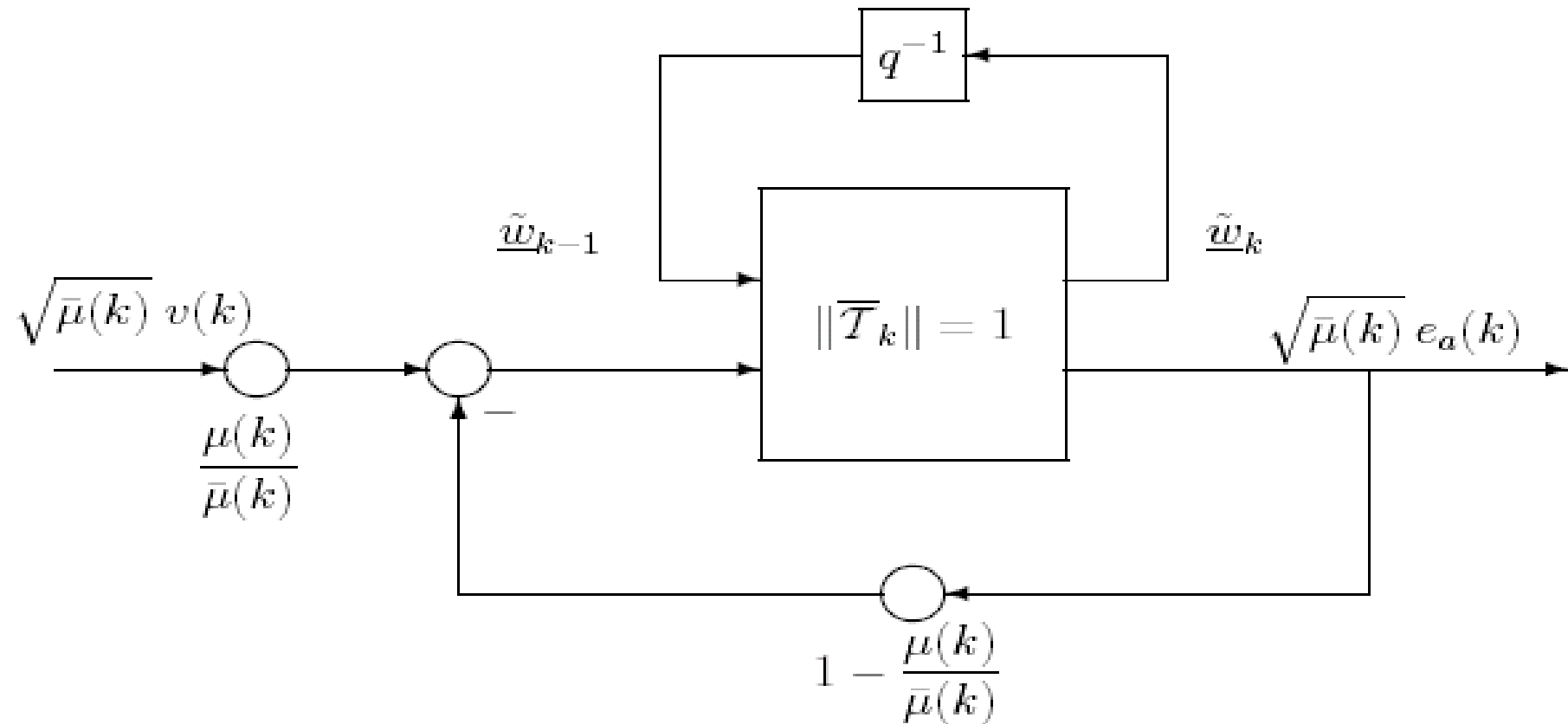
$$\begin{aligned}\hat{\underline{w}}_k &= \hat{\underline{w}}_{k-1} + \mu_k \underline{x}_k^* [e_{a,k} + v_k] \\ &= \hat{\underline{w}}_{k-1} + \bar{\mu}_k \underline{x}_k^* [e_{a,k} + \bar{v}_k]\end{aligned}$$

- using

$$\begin{aligned}-\bar{v}_k &= e_{a,k} - \frac{\mu_k}{\bar{\mu}_k} [e_{a,k} + v_k] \\ &= \left[1 - \frac{\mu_k}{\bar{\mu}_k} \right] e_{a,k} - \frac{\mu_k}{\bar{\mu}_k} v_k\end{aligned}$$



Feedback



Convergence of Feedback Scheme

- Consider again
$$\sum_{k=1}^N \bar{\mu}_k |e_{a,k}|^2 \leq \|\underline{w} - \hat{w}_0\|_2^2 + \sum_{k=1}^N \bar{\mu}_k |\bar{v}_k|^2$$

$$\bar{v}_k \stackrel{\Delta}{=} \left[\frac{\mu_k}{\bar{\mu}_k} \right] v_k + \left[\frac{\mu_k}{\bar{\mu}_k} - 1 \right] e_{a,k}$$

- Using short notations:

$$\delta(N) \stackrel{\Delta}{=} \max_{1 \leq k \leq N} \left| 1 - \frac{\mu(k)}{\bar{\mu}(k)} \right|$$

$$\gamma(N) \stackrel{\Delta}{=} \max_{1 \leq k \leq N} \frac{\mu(k)}{\bar{\mu}(k)}$$



Convergence of Feedback Scheme

$$\sqrt{\sum_{k=1}^N \bar{\mu}_k |e_{a,k}|^2} \leq \sqrt{\|\underline{\tilde{w}}_0\|_2^2} + \sqrt{\sum_{k=1}^N \bar{\mu}_k |\bar{v}_k|^2}$$

Use triangular inequality

$$\leq \sqrt{\|\underline{\tilde{w}}_0\|_2^2} + \sqrt{\sum_{k=1}^N \bar{\mu}_k \left[\frac{\mu_k}{\bar{\mu}_k} \right]^2 |v_k|^2} + \sqrt{\sum_{k=1}^N \bar{\mu}_k \left[1 - \frac{\mu_k}{\bar{\mu}_k} \right]^2 |e_{a,k}|^2}$$

$$\leq \|\underline{\tilde{w}}_0\|_2 + \max_k \left| \frac{\mu_k}{\bar{\mu}_k} \right| \sqrt{\sum_{k=1}^N \bar{\mu}_k |v_k|^2} + \max_k \left| 1 - \frac{\mu_k}{\bar{\mu}_k} \right| \sqrt{\sum_{k=1}^N \bar{\mu}_k |e_{a,k}|^2}$$

$$= \|\underline{\tilde{w}}_0\|_2 + \gamma_N \sqrt{\sum_{k=1}^N \bar{\mu}_k |v_k|^2} + \delta_N \sqrt{\sum_{k=1}^N \bar{\mu}_k |e_{a,k}|^2}$$

$$= \frac{1}{1 - \delta_N} \left[\|\underline{\tilde{w}}_0\|_2 + \gamma_N \sqrt{\sum_{k=1}^N \bar{\mu}_k |v_k|^2} \right]$$



Convergence of Feedback Scheme

- **Theorem 5.4 (Extended Convergence of the Gradient Type Algorithm)**

$$\text{If } \mu_k \|\underline{x}_k\|^2 < 2, \|\underline{\tilde{w}}_0\| < \infty, \text{ and } \sum_{k=1}^{\infty} \bar{\mu}_k |v_k|^2 < \infty$$

then it follows that: $\sqrt{\bar{\mu}_k} e_{a,k} \rightarrow 0$.

If, furthermore, $\sqrt{\mu_k} \underline{x}_k$ is persistent exciting, then we also have:

$$\underline{\hat{w}}_k \rightarrow \underline{w}$$



Convergence of Feedback Scheme

- Also consider the energy flow in such system. Since the forward path has no loss of energy, all energy going in must come out. The energy part in the parameter error vector is fed back. Only here in the feedback path energy can be lost. The more energy is lost, the faster is the adaptation. Fastest adaptation is thus obtained for

$$\mu_k = \overline{\mu}_k$$



Nonlinear, Memoryless Filter in the Estimation Path

- Consider now $f[y_k]$ instead of y_k
- Perceptron Learning Algorithm



Perceptron Learning Algorithm

- Consider two sets S_1 and S_2 of M -dimensional real-valued vectors \underline{x}_k

$$S_1 = \{ \underline{x}_k \in R^M \mid \underline{x}_k \text{ has property } A \}$$

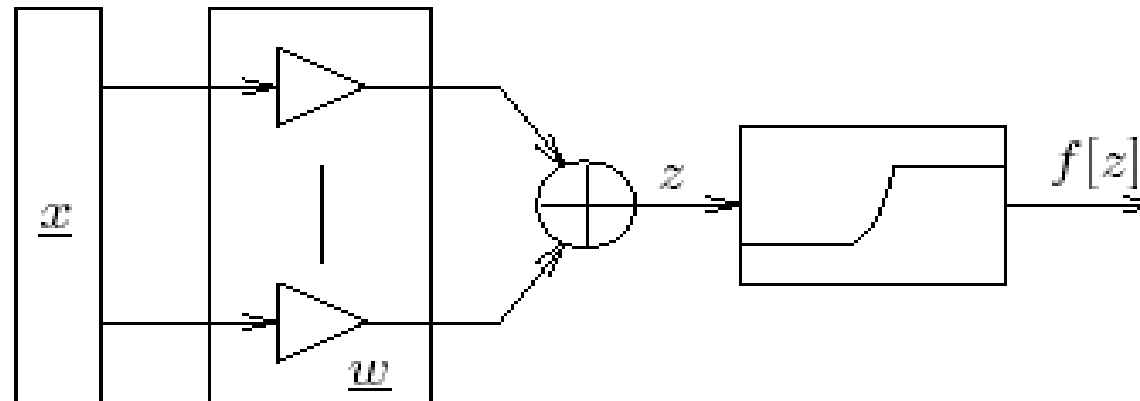
$$S_2 = \{ \underline{x}_k \in R^M \mid \underline{x}_k \text{ has property } \overline{A} \}$$

- If both sets are linearly separable, then they describe two classes.



Perceptron Learning Algorithm

- Model of linear perceptron



Perceptron Learning Algorithm

- Activity function is in general a sigmoid-function

- In particular:
$$f_{\beta}[z] = \frac{1}{1 + e^{-\beta z}}; \quad \beta > 0$$

$$f_{\infty}[z] = \frac{1 + \text{sgn}[z]}{2} = \begin{cases} 0 & ; \text{if } z < 0 \\ 1/2 & ; \text{if } z = 0 \\ 1 & ; \text{if } z > 0 \end{cases}$$



Perceptron Learning Algorithm

- Consider input vectors \underline{x}_k with their corresponding output values y_k

$$y_k = f[\underline{x}_k^T \underline{w}]$$

- In supervised learning, the data pairs $\{\underline{x}_k, y_k\}$ are offered to the synapse in order to find \underline{w} .



Perceptron Learning Algorithm

- The most well-known algorithm with such property is the PLA

$$\underline{\hat{w}}_k = \underline{\hat{w}}_{k-1} + \mu \underline{x}_k \left(y_k - f \left[\underline{x}_k^T \underline{\hat{w}}_{k-1} \right] \right)$$

- Consider again the general form including additive noise

$$d_k = f \left[\underline{x}_k^T \underline{w} \right] + v_k = y_k + v_k$$

- This leads to:

$$\underline{\hat{w}}_k = \underline{\hat{w}}_{k-1} + \mu_k \underline{x}_k \left(d_k - f \left[\underline{x}_k^T \underline{\hat{w}}_{k-1} \right] \right)$$



Perceptron Learning Algorithm

- Note that this update

$$\underline{\hat{w}}_k = \underline{\hat{w}}_{k-1} + \mu \underline{x}_k \left(y_k - f \left[\underline{x}_k^T \underline{\hat{w}}_{k-1} \right] \right)$$

- is a somewhat simplified form of

$$\underline{\hat{w}}_k = \underline{\hat{w}}_{k-1} + \tilde{\mu} \underline{x}_k \left(y_k - f \left[\underline{x}_k^T \underline{\hat{w}}_{k-1} \right] \right) f' \left[\underline{x}_k^T \underline{\hat{w}}_{k-1} \right]$$

- Where we incorporated the derivative into the stepsize



Perceptron Learning Algorithm

- Because of the mean value theorem we have

$$f[\underline{x}_k^T \underline{w}] - f[\underline{x}_k^T \hat{\underline{w}}_{k-1}] = f'[\eta_k] e_{a,k}$$

- For
- With this, the PLA can be treated as the previous algorithms.

$$\underline{x}_k^T \underline{w} < \eta_k < \underline{x}_k^T \hat{\underline{w}}_{k-1}$$



Perceptron Learning Algorithm

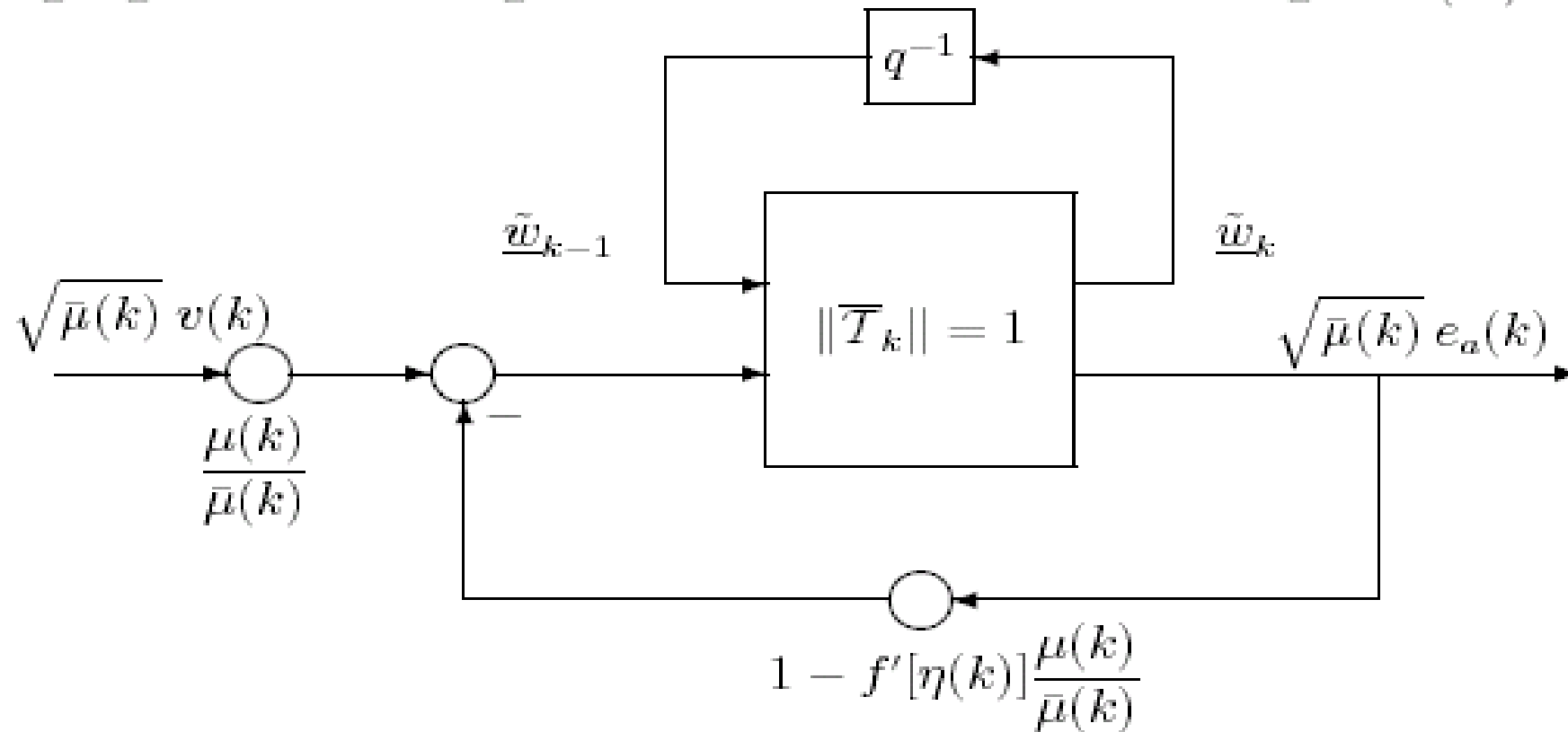
- We find again the stability condition

$$\delta_N = \max_{1 \leq k \leq N} \left| 1 - f'[\eta_k] \frac{\mu_k}{\bar{\mu}_k} \right|$$

- Together with a feedback structure



Perceptron Learning Algorithm



PLA

- Different to the LMS algorithm, the PLA includes a derivative of the nonlinear function in its feedback path
- As we (mostly) do not know in advance at which point we have to take the derivative, we take its largest value as upper bound
- Sigmoid function: max derivative is $\frac{1}{4} \beta$
- Tanh function: max derivative is $\frac{1}{2} \beta$
- Leads to conservative results but ensures stability!



PLA

- Recall that we have incorporated the derivative into the step-size
- If we would not have done so, the derivative would appear twice in the feedback loop
 - However, not exactly the same term as one goes with

$$\underline{x}_k^T \hat{\underline{w}}_{k-1}$$

- While the other goes with

$$\underline{x}_k^T \underline{w} < \eta_k < \underline{x}_k^T \hat{\underline{w}}_{k-1}$$

- We will see in the backpropagation algorithm, the update will change...



Linear and Non-Linear Filter in Reference Path

- Recurrent Neuronal Networks (RNN) are very suitable for classification
- The structure proposed by Narendra and Parthasarathy is often treated in literature since it is suitable for calculation.



Linear and Non-Linear Filter in Reference Path

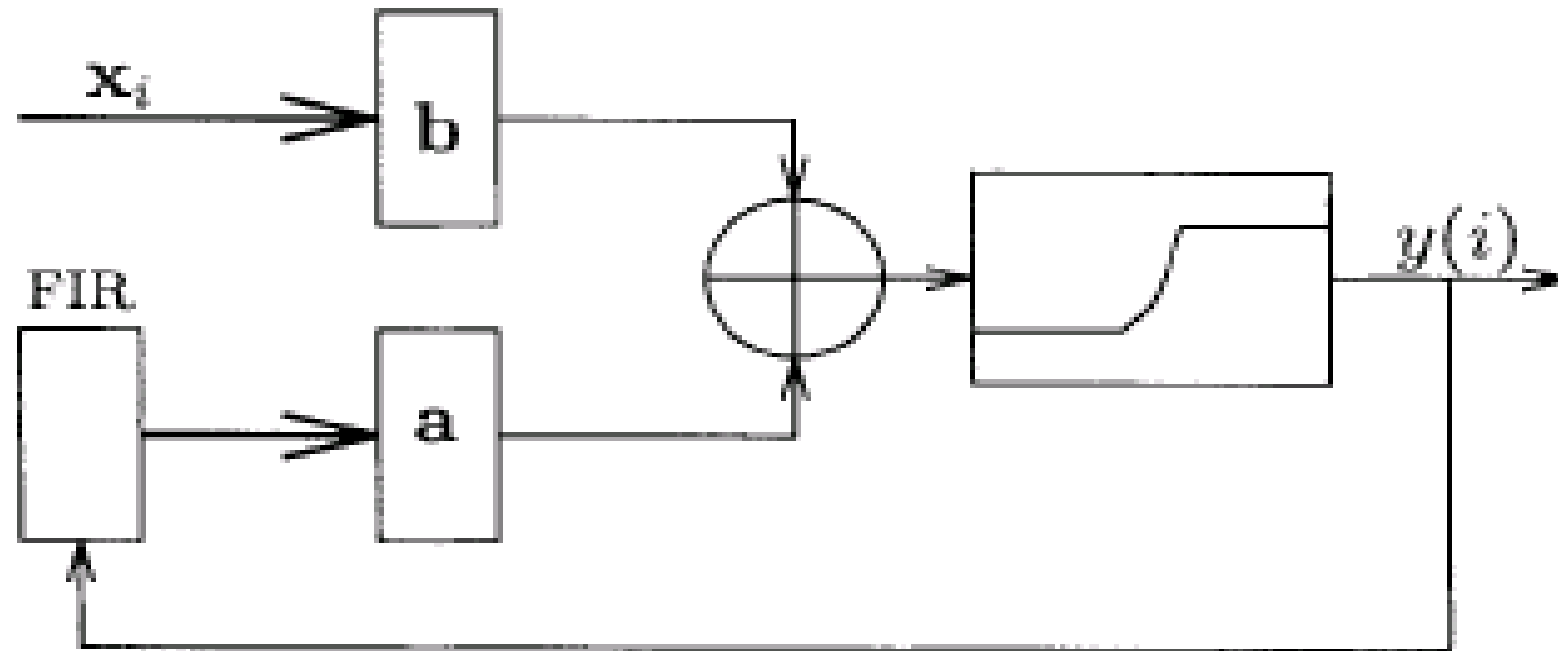


Fig. 6. Narendra and Parthasarathy's dynamic network.



Linear and Non-Linear Filter in Reference Path

- The update equation can be written as:

- With the values
$$\underline{\hat{w}}_k = \underline{\hat{w}}_{k-1} + \mu_k \underline{u}_k \left[\frac{1}{f'^{-1}[\eta_k] - A(q^{-1})} [e_{a,k}] + v_k \right]$$

- We find:

$$\lambda = \min_k \{f'^{-1}[\eta_k]\}; \quad \zeta^{-1} = \max_k \{\bar{\mu}_k^{-1}\}$$

$$\text{Real}\{1 - A(e^{j\Omega})\} < \delta, \quad \Omega \in [0, 2\pi]$$

$$\mu < 2\zeta(\lambda + \delta - 1)$$



Linear and Non-Linear Filter in Reference Path

- There is a direct relation to the slope of the derivative of the sigmoid function:

$$\beta < \frac{8}{\frac{\mu}{\zeta} + 2(1 - \delta)}$$

- Note that although nonlinearity and a linear filter occur in the reference path, it is equivalently a system with nonlinear filter in the error path!



Literature

- M. Rupp, A.H. Sayed:
"*Supervised Learning of Perceptron and Output Feedback Dynamic Networks: A Feedback Analysis via the Small Gain Theorem*";
IEEE Transactions on Neural Networks, **Vol. 8** (1997), no. 3; S. 612 - 622.

