

LECTURES ON  
NEURAL NETWORKS

BY

PROF. BERNARD WIDROW

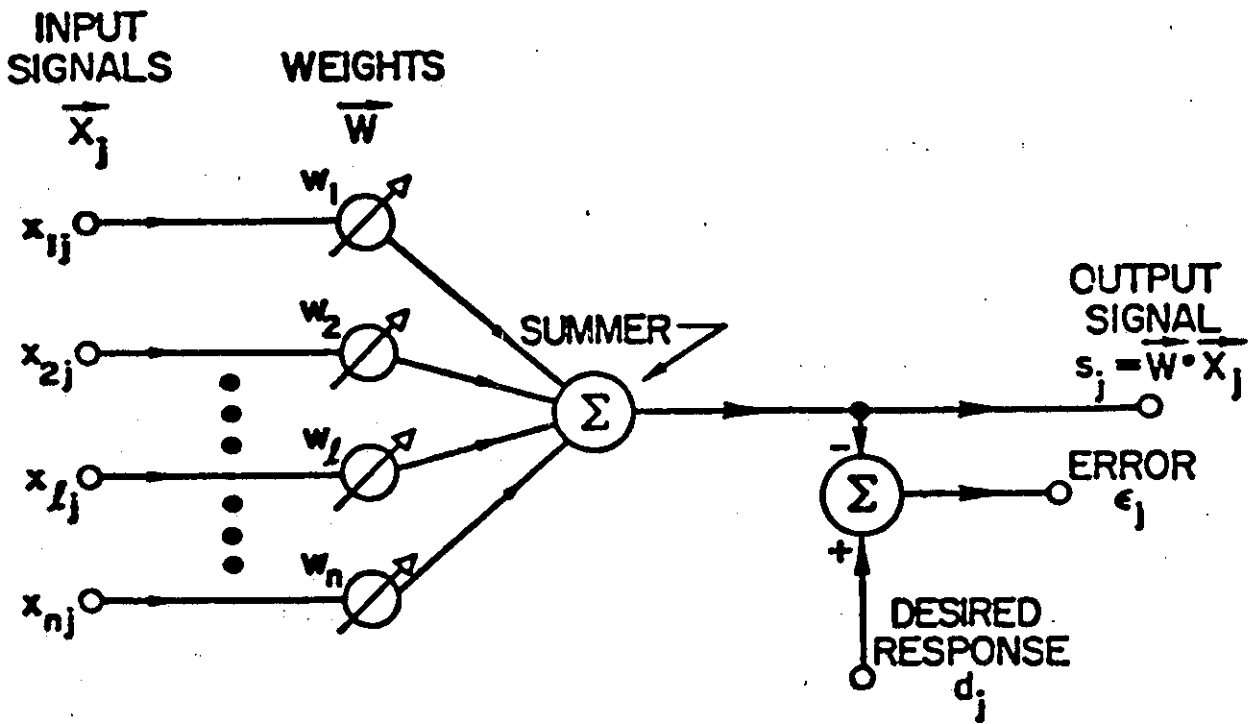
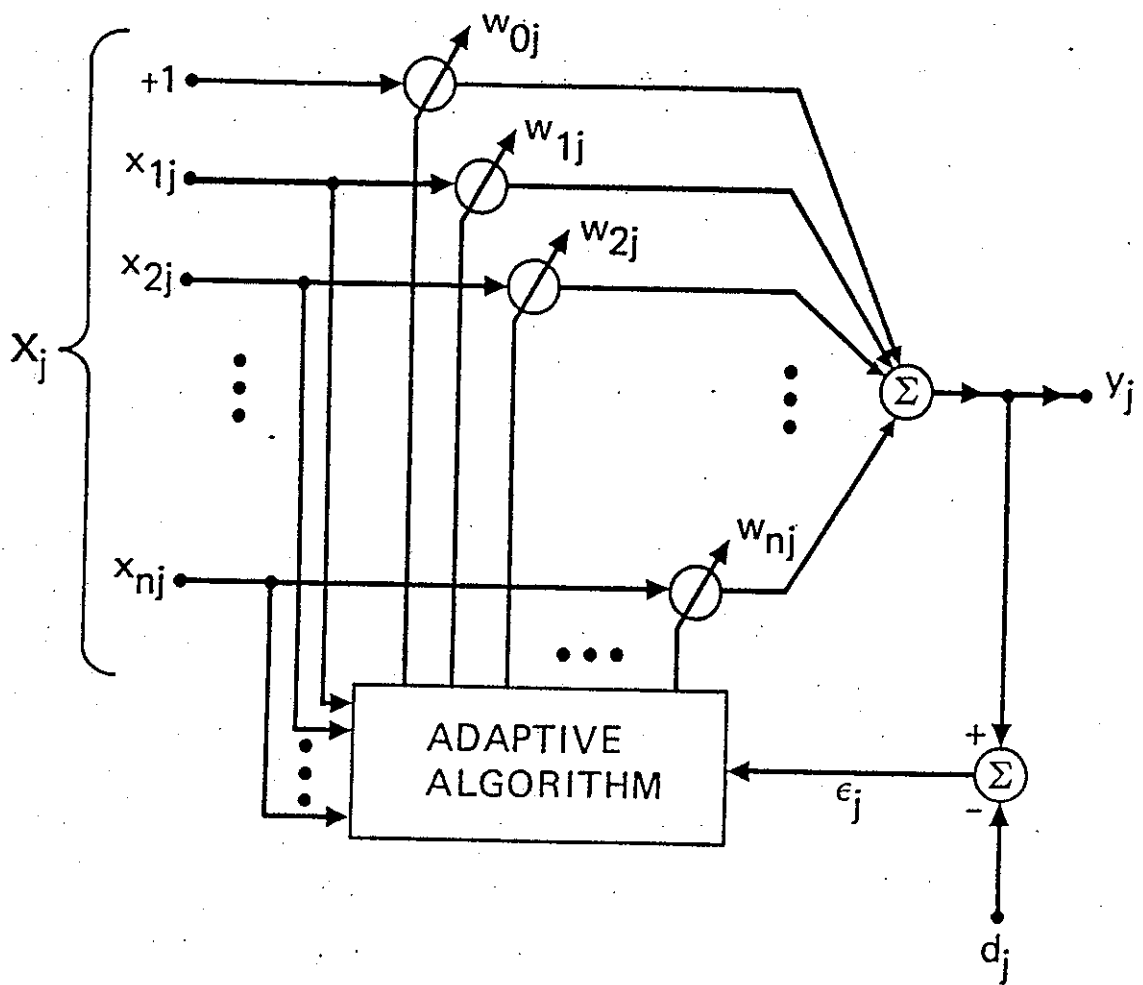
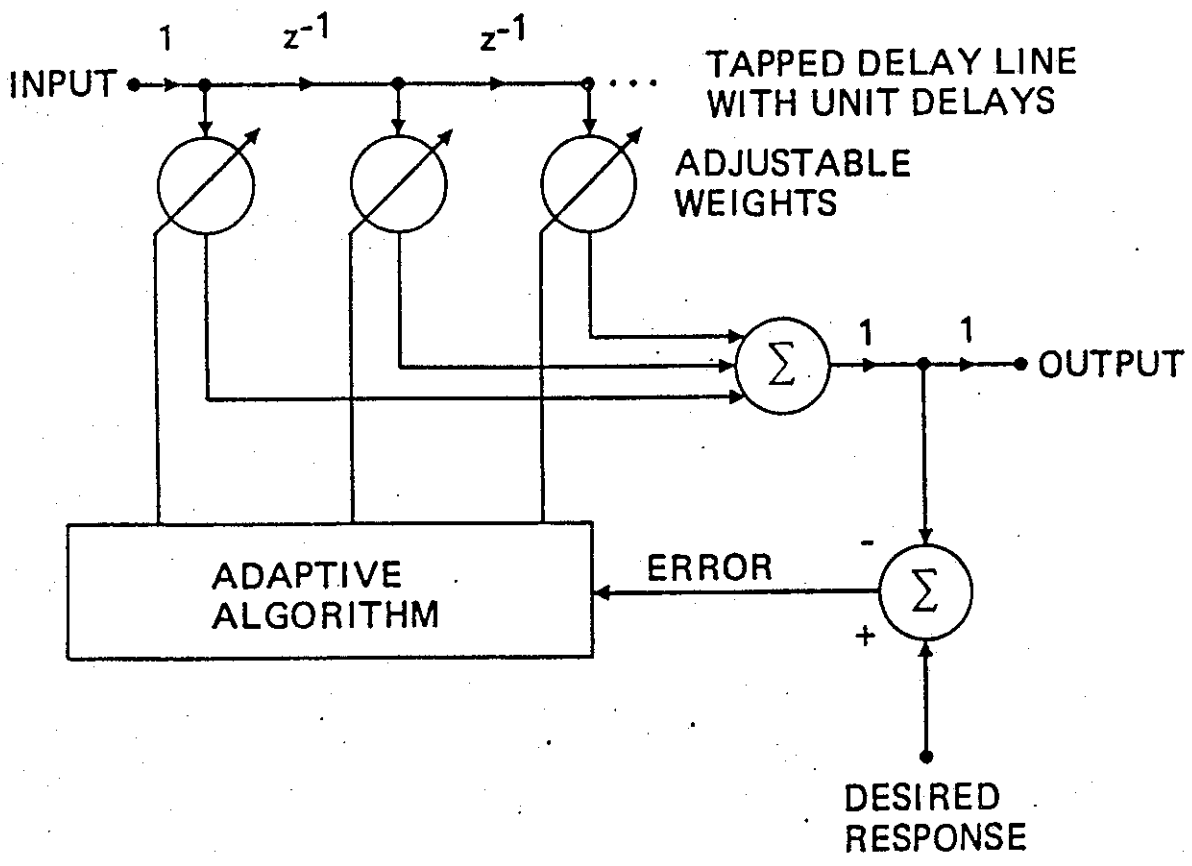


Fig. 2. Adaptive linear combinatorial system.





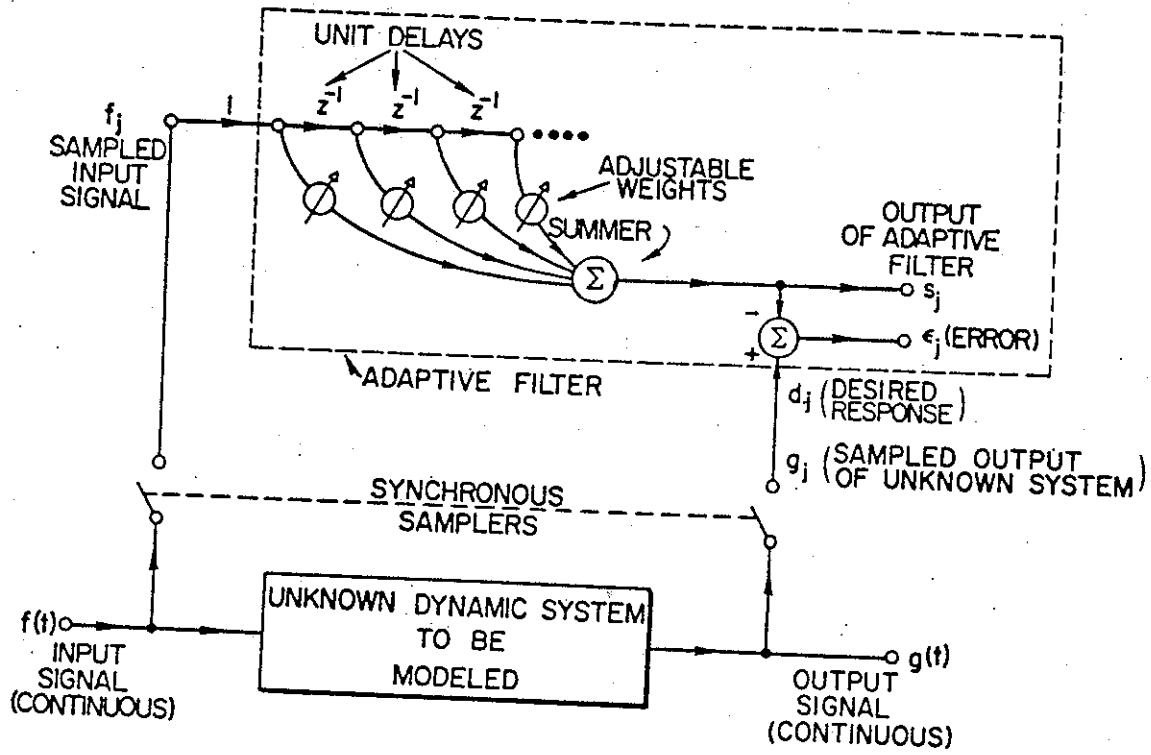
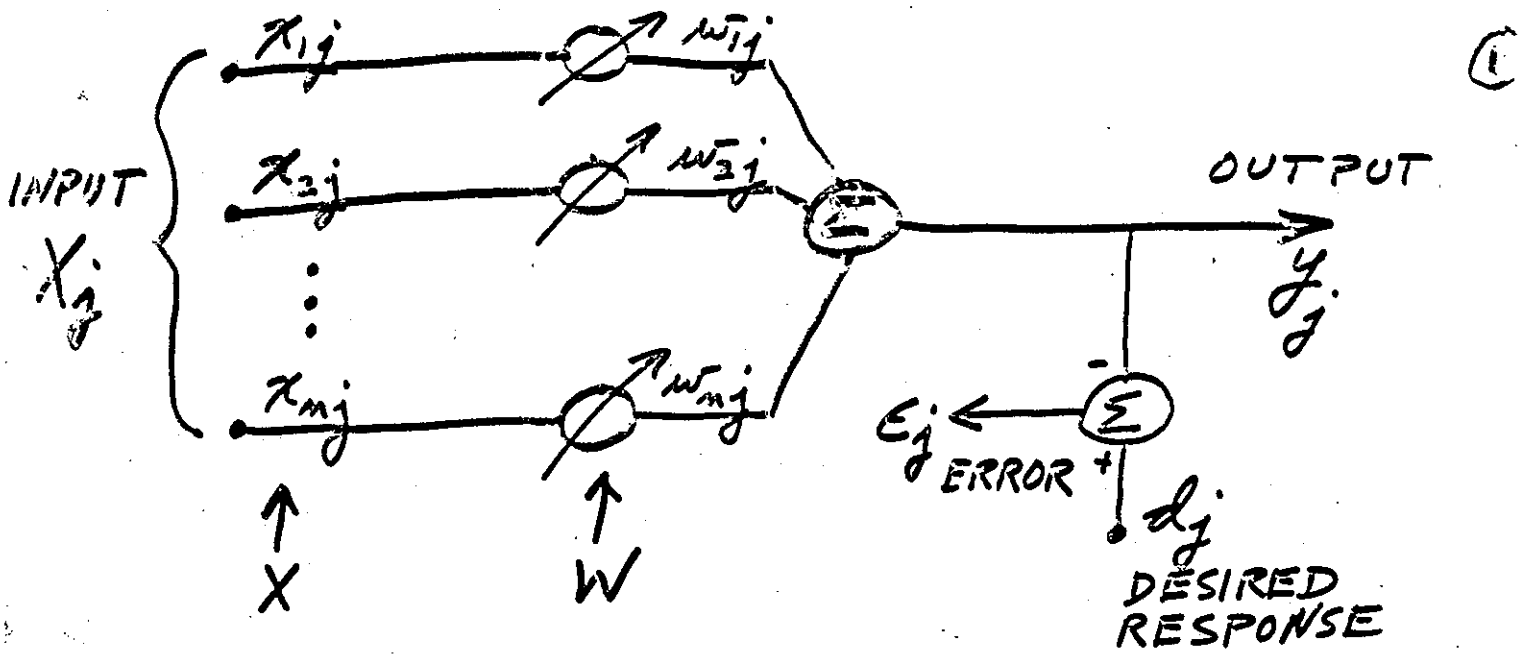


Fig. 1. Modeling an unknown system by a discrete adaptive filter.



$$y_j = X_j^T W = W^T X_j$$

$$e_j = d_j - X_j^T W$$

$$e_j^2 = d_j^2 - 2d_j X_j^T W + W^T X_j X_j^T W$$

$$MSE = \sum e_j^2 \triangleq E[e_j^2] = E[d_j^2] - 2P^T W + W^T R W$$

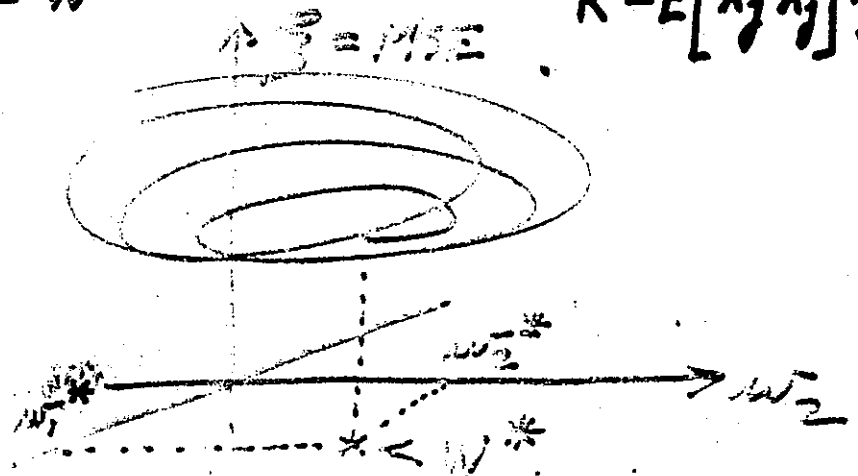
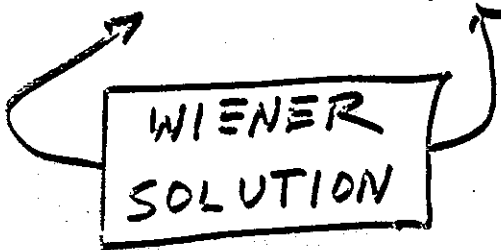
$$\nabla = \frac{\partial \sum}{\partial W} = -2P + 2RW$$

$$P \triangleq E\left\{ \begin{matrix} \leftarrow 1 \rightarrow \\ d_j X_j \end{matrix} \right\} \begin{matrix} \uparrow \\ \leftarrow n \rightarrow \\ \downarrow \end{matrix}$$

$$R \triangleq E\left[ \begin{matrix} \leftarrow n \rightarrow \\ X_j X_j^T \\ \downarrow \end{matrix} \right]$$

when  $\nabla = 0$ ,  $W = W^*$

$$\therefore W^* = R^{-1} P$$



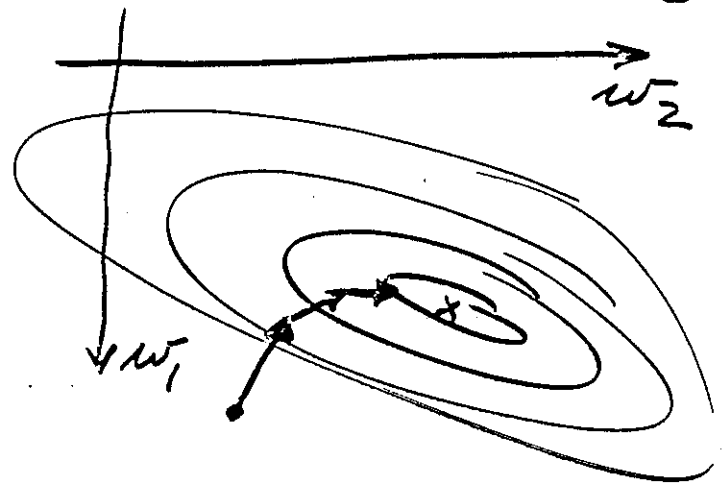
STEEPEST DESCENT:

(2)

$$W_{j+1} = W_j - \mu \nabla_j$$

TRUE GRADIENT:

$$\nabla = \frac{\partial \mathcal{E}}{\partial W} = -2P + 2RW$$



A GRADIENT ESTIMATE:

$$\hat{\nabla} = \frac{\partial \epsilon_j^2}{\partial W} \equiv \frac{\partial E[\epsilon_j^2]}{\partial W} \triangleq \nabla$$

now,  $\epsilon_j = d_j - X_j^T W$

$$\hat{\nabla} = 2\epsilon_j \frac{\partial \epsilon_j}{\partial W} = 2\epsilon_j (-X_j) = -2\epsilon_j X_j$$

$$E[\hat{\nabla}] = -2E[\epsilon_j X_j] = -2E[d_j X_j - X_j X_j^T W] = -2P + 2RW = \nabla$$

$\therefore$  grad. estimate is unbiased!

LMS ALGORITHM (Least mean square error)  
OF WIDROW and HOFF (1959):

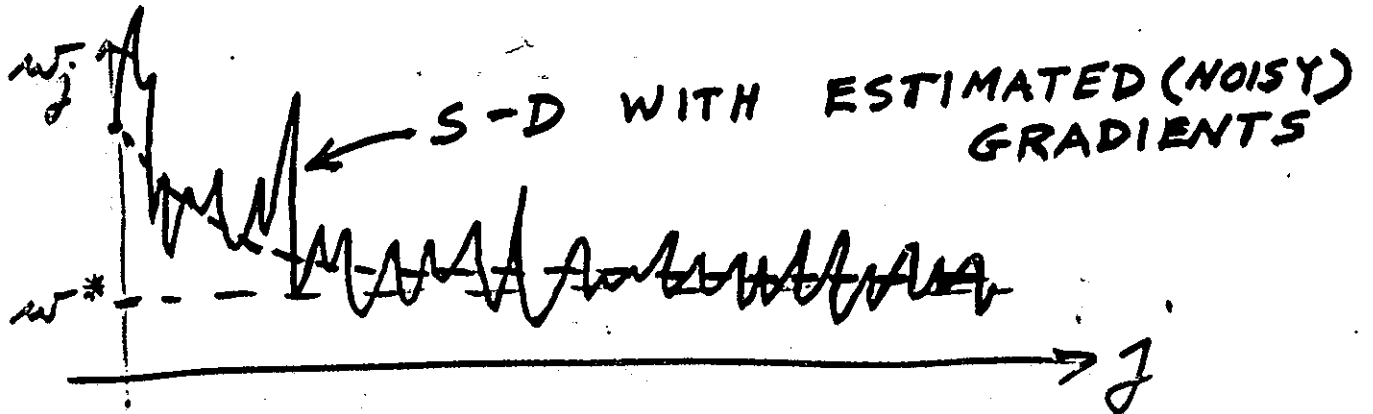
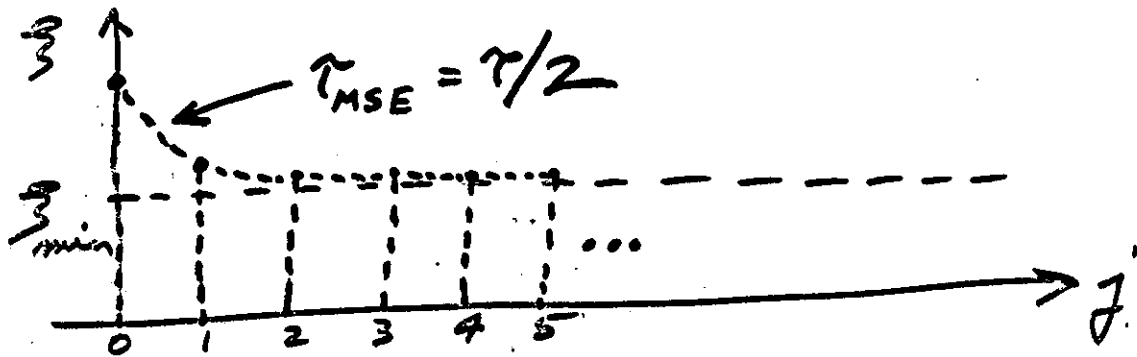
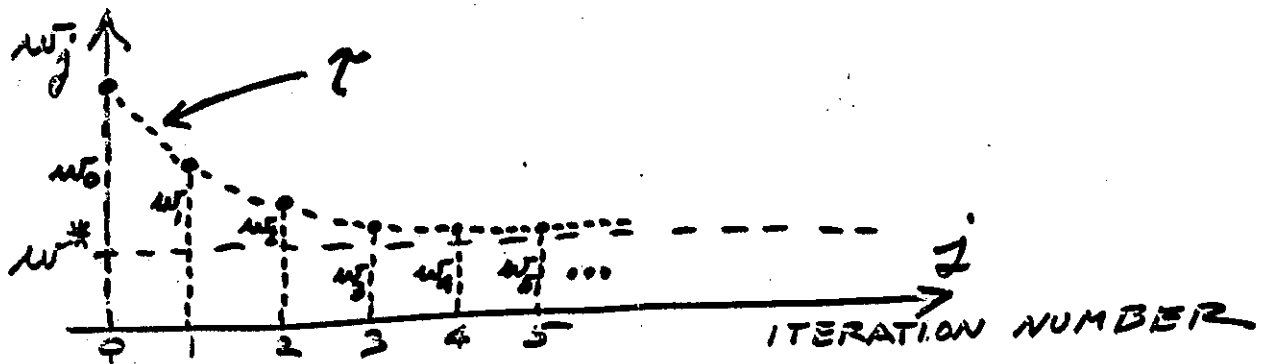
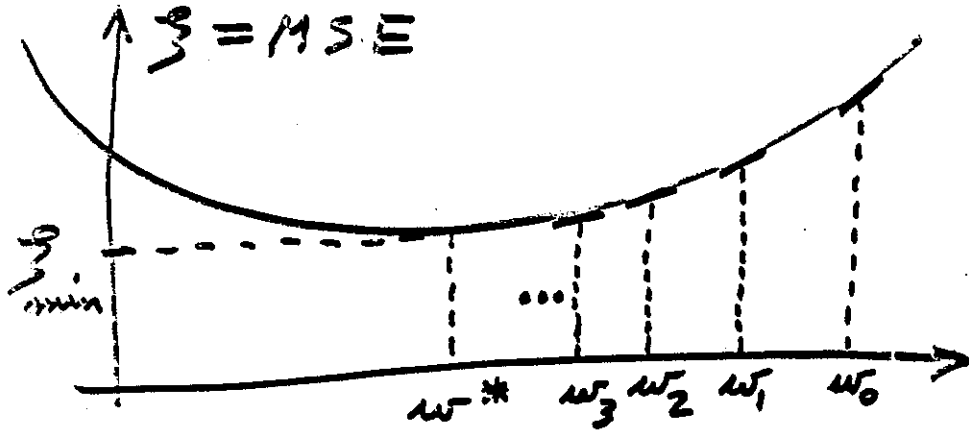
$$W_{j+1} = W_j - \mu \hat{\nabla}_j \leftarrow \text{S-D using estimated gradient}$$

$$W_{j+1} = W_j + 2\mu \epsilon_j X_j$$

and  $\epsilon_j = d_j - X_j^T W_j$

 $\leftarrow$  LMS

3





CONDITION FOR CONVERGENCE (STABILITY) OF THE MEAN OF THE WEIGHT VECTOR: (5)

$$\lim_{j \rightarrow \infty} [I - 2\mu \Lambda]^j = 0, \text{ i.e.}$$

$$\frac{1}{\lambda_{\max}} > \mu > 0$$

SUFFICIENT (BUT NOT NECESSARY), EASY TO APPLY, CONVERGENCE OF THE MEAN:

$$\frac{1}{\text{TRACER}} > \mu > 0$$

NOTE: CONVERGENCE OF THE MEAN OF THE WEIGHT VECTOR DOES NOT GUARANTEE CONVERGENCE OF ITS VARIANCE OF OTHER MOMENT.

TIME CONSTANTS OF LMS ALGORITHM }  $\tau_p = \frac{1}{2\mu \lambda_p}$

$\tau_{pMSE} = \frac{1}{4\mu \lambda_p}$

MISADJUSTMENT:

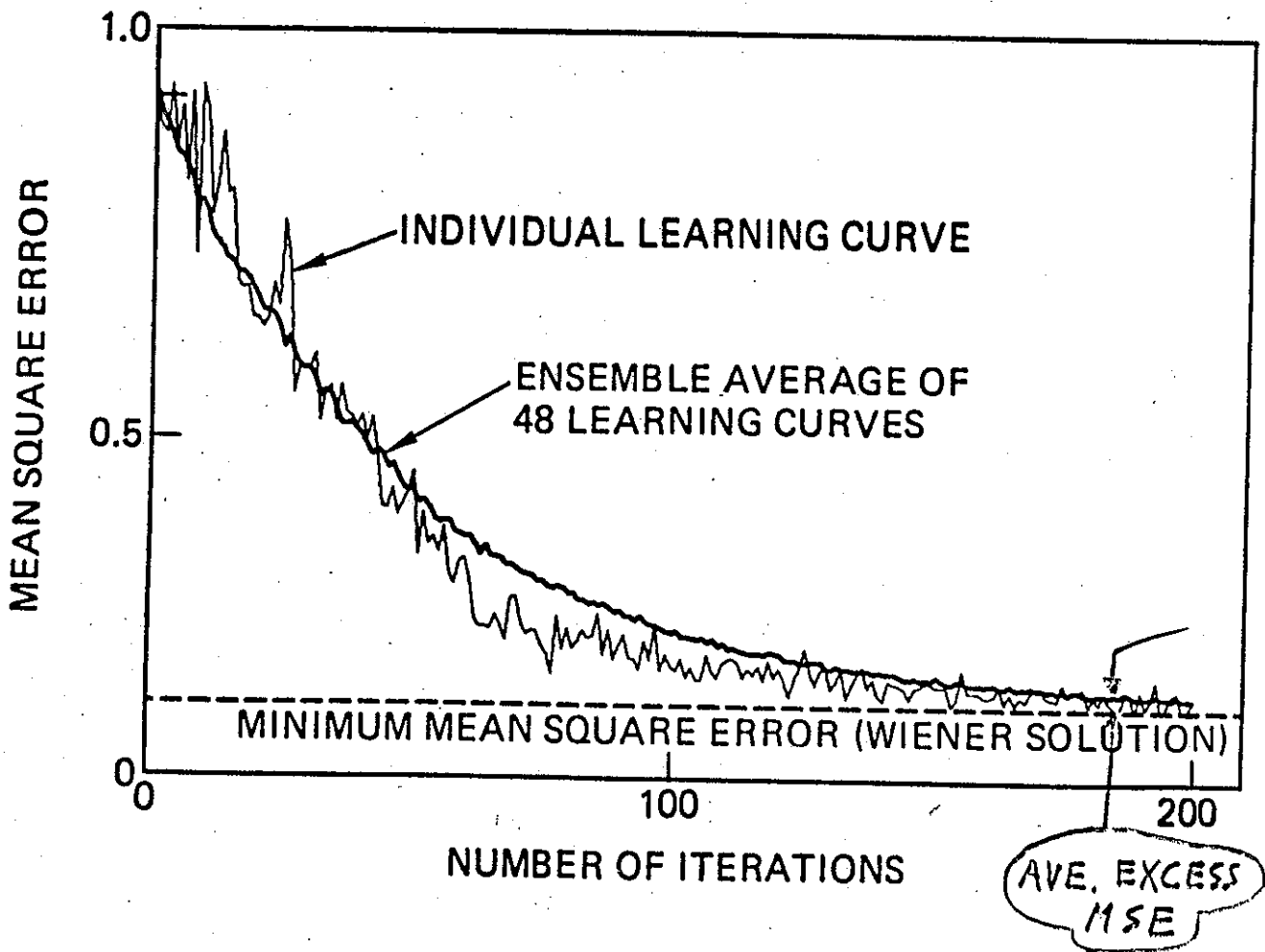
$$M \triangleq \frac{\text{AVE EXCESS MSE}}{\text{MINIMUM MSE}} = \mu \text{ TRACER} = \frac{1}{2} \sum_{p=1}^n \frac{1}{\tau_p} = \frac{n}{2} \left( \frac{1}{\tau_p} \right)_{\text{AVE}}$$

$$M = \frac{n}{4} \left( \frac{1}{\tau_{pMSE}} \right)_{\text{AVE}} = \frac{n}{4} \frac{1}{\tau_{MSE}} \leftarrow \text{WHEN ALL } \lambda\text{'S ARE EQUAL}$$

EXAMPLE: Let  $M = 10\%$  be an "OK" misadj, and let  $n = 10$  weights. The question is, how big does  $\tau_{MSE}$  need to be, and what is the settling time of the adaptive process?

$$M = 0.1 = \frac{10}{4} \frac{1}{\tau_{MSE}} \quad \therefore \tau_{MSE} = 25 \text{ iterations}$$

Settling time  $\approx 100$  iterations = 10X filter length.

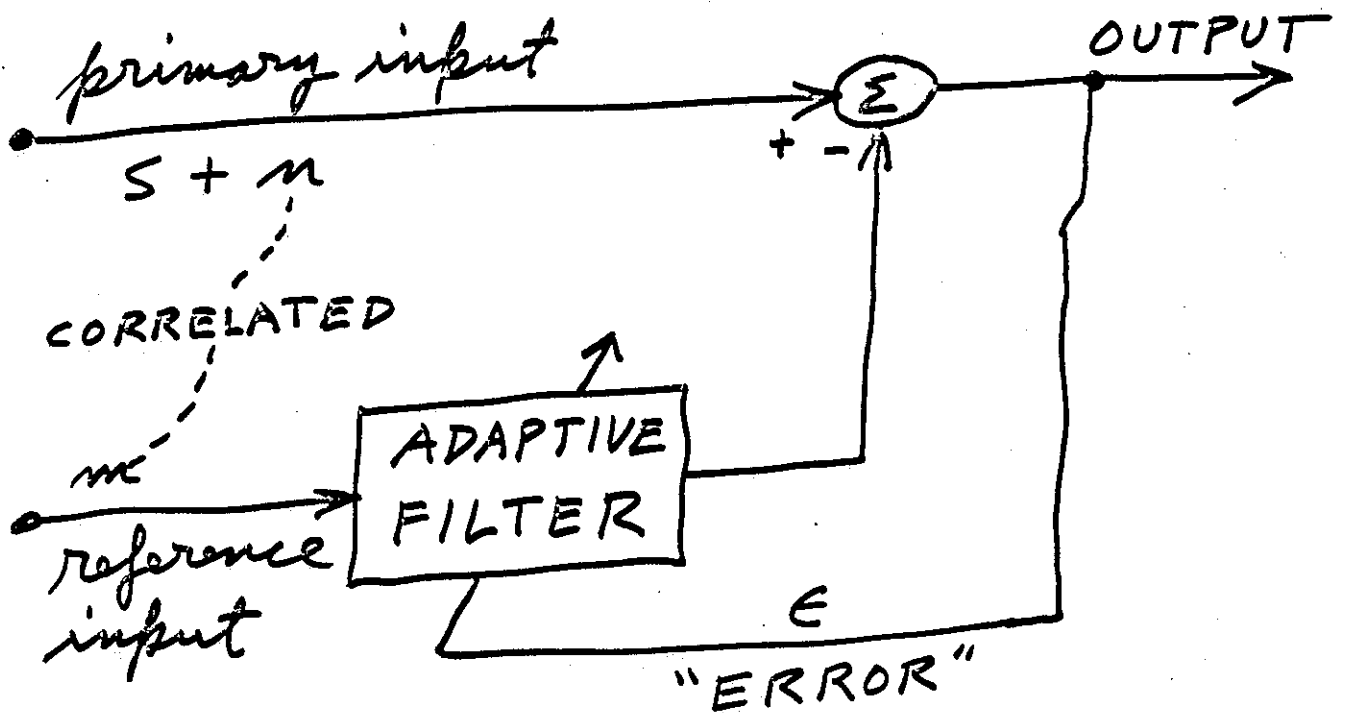


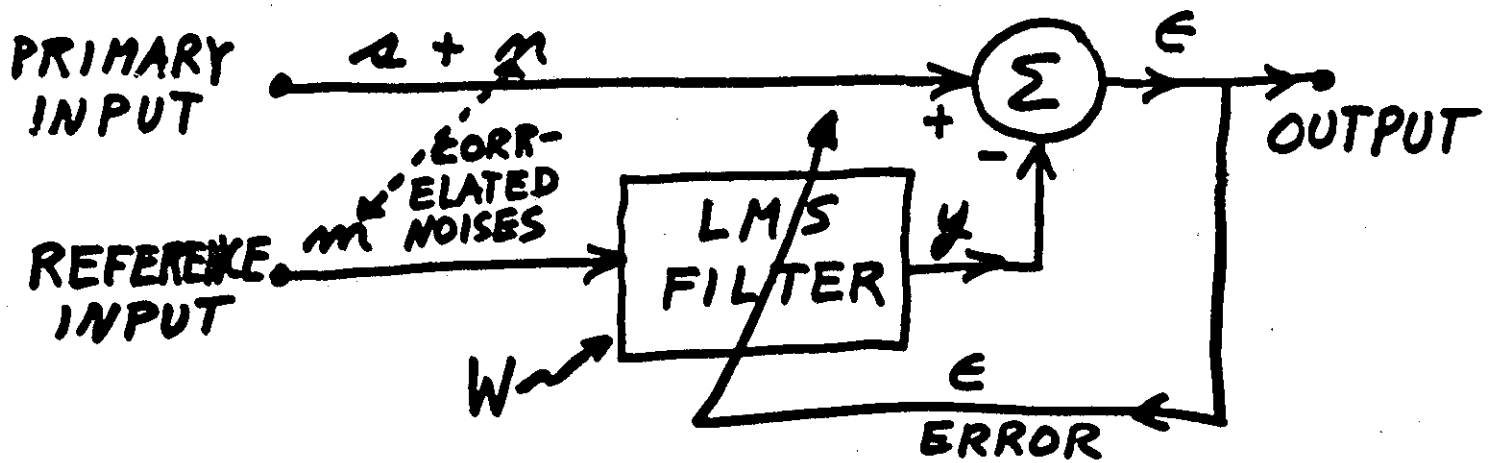
# CONVENTIONAL FILTERING



---

# ADAPTIVE NOISE CANCELLING





$$E = a + n - y$$

$$E^2 = [a + (n - y)]^2 = a^2 + (n - y)^2 + 2a(n - y)$$

$$E[E^2] = E[a^2] + E[(n - y)^2] + 2E[a(n - y)]$$

$$E[E^2] = E[a^2] + E[(n - y)^2].$$

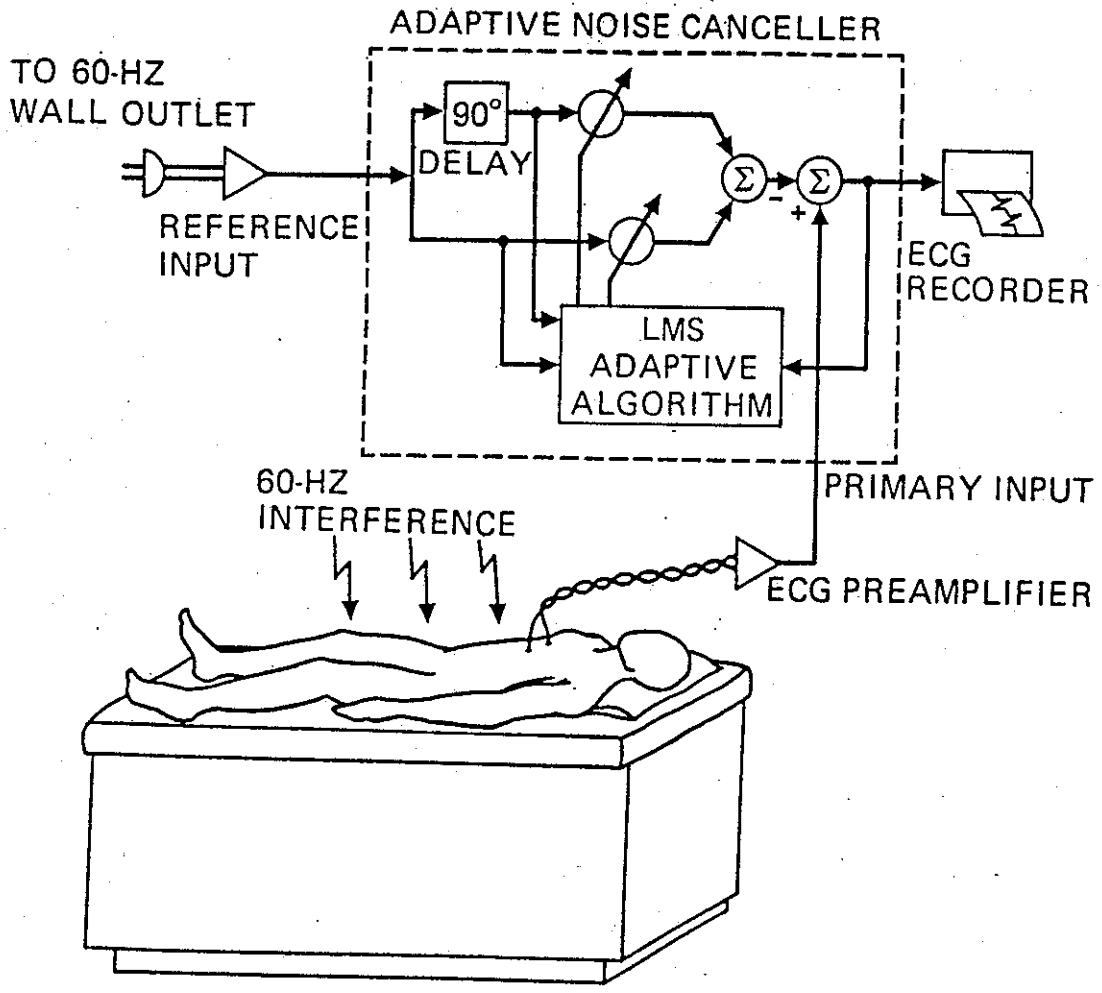
$$\min_W E[E^2] = E[a^2] + \min_W [E(n - y)^2], \text{ THUS}$$

(1) MINIMIZING  $E[E^2]$  MINIMIZES  $E[(n - y)^2]$ ,

(2) CAUSES  $y$  TO BE BEST LEAST SQUARES ESTIMATE OF  $m$ ,

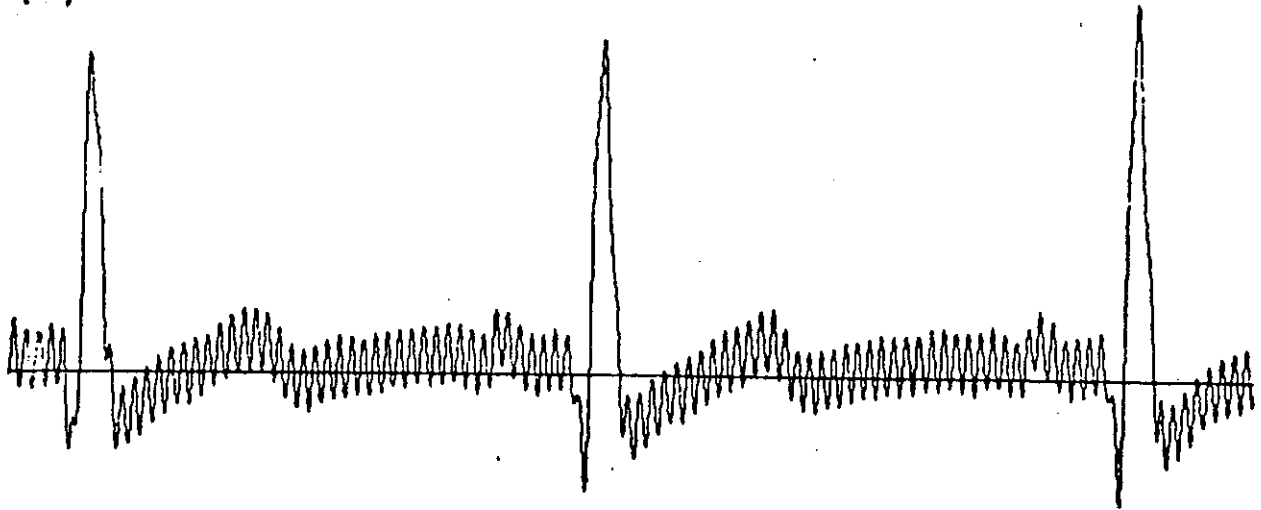
(3) " " " " " " " " OF  $a$ ,

(4) MAXIMIZES OUTPUT SNR WITHOUT SIGNAL DISTORTION

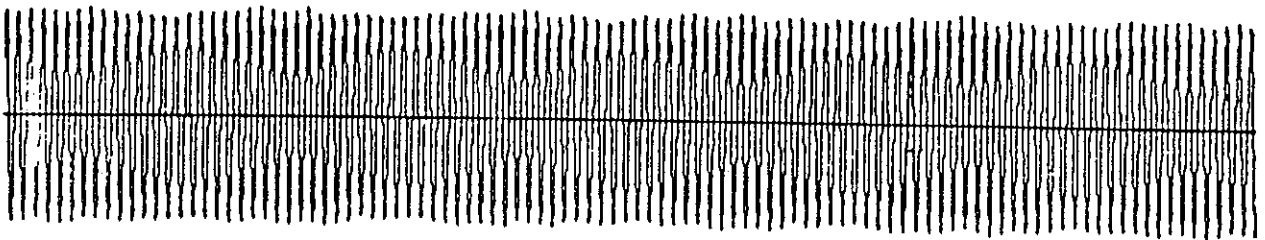


CANCELLING 60HZ INTERFERENCE IN THE EKG

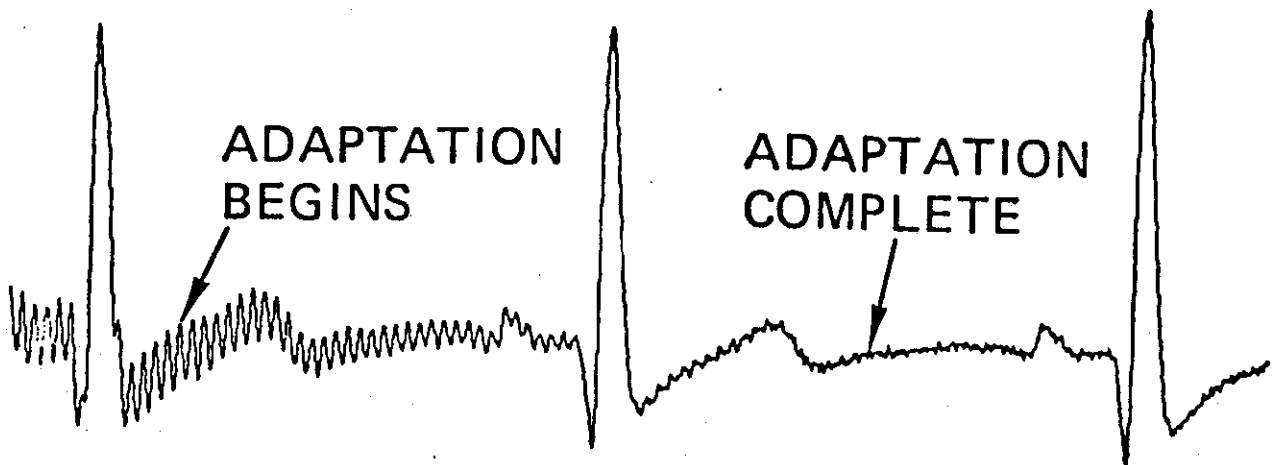
(a)



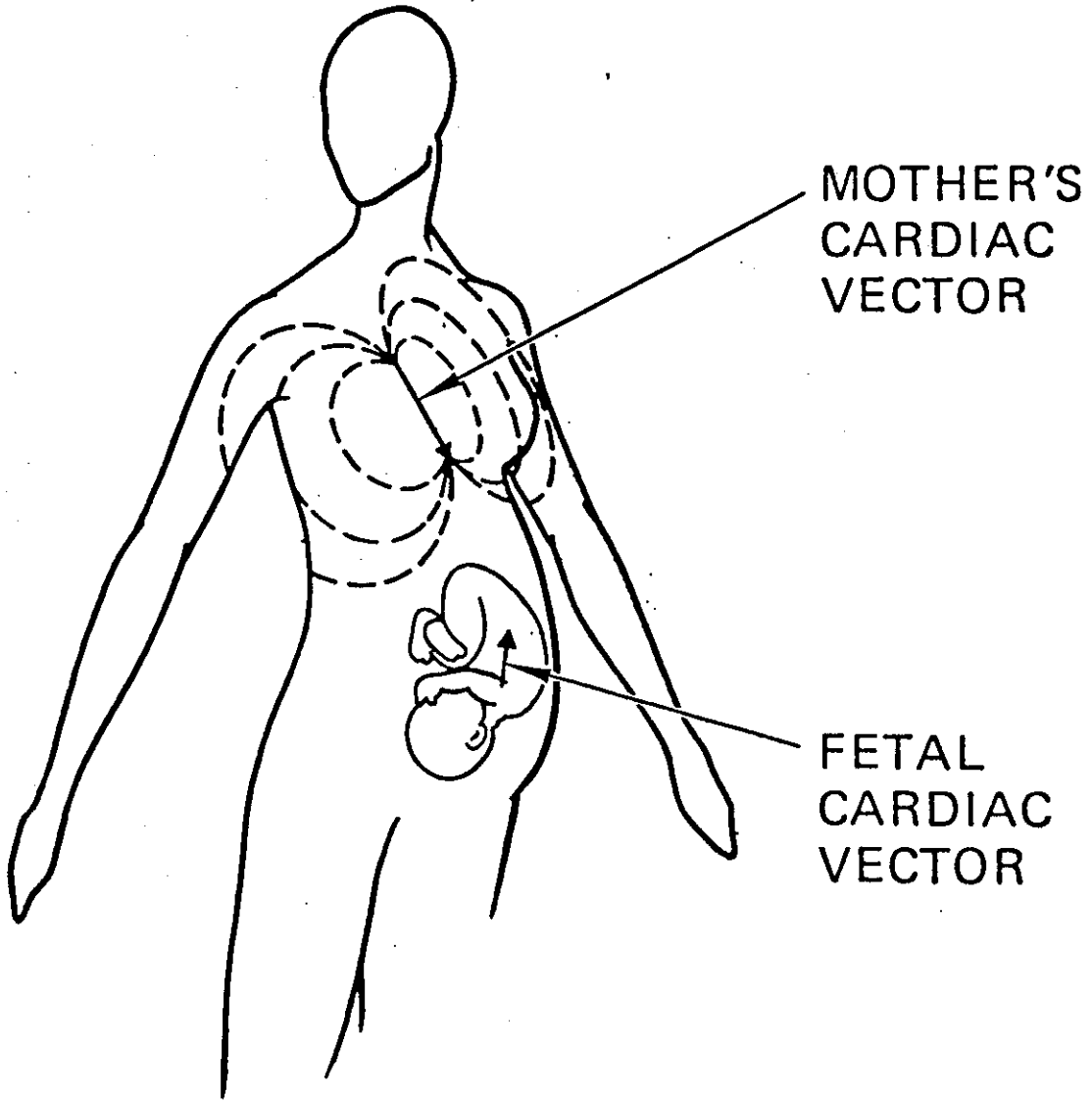
(b)



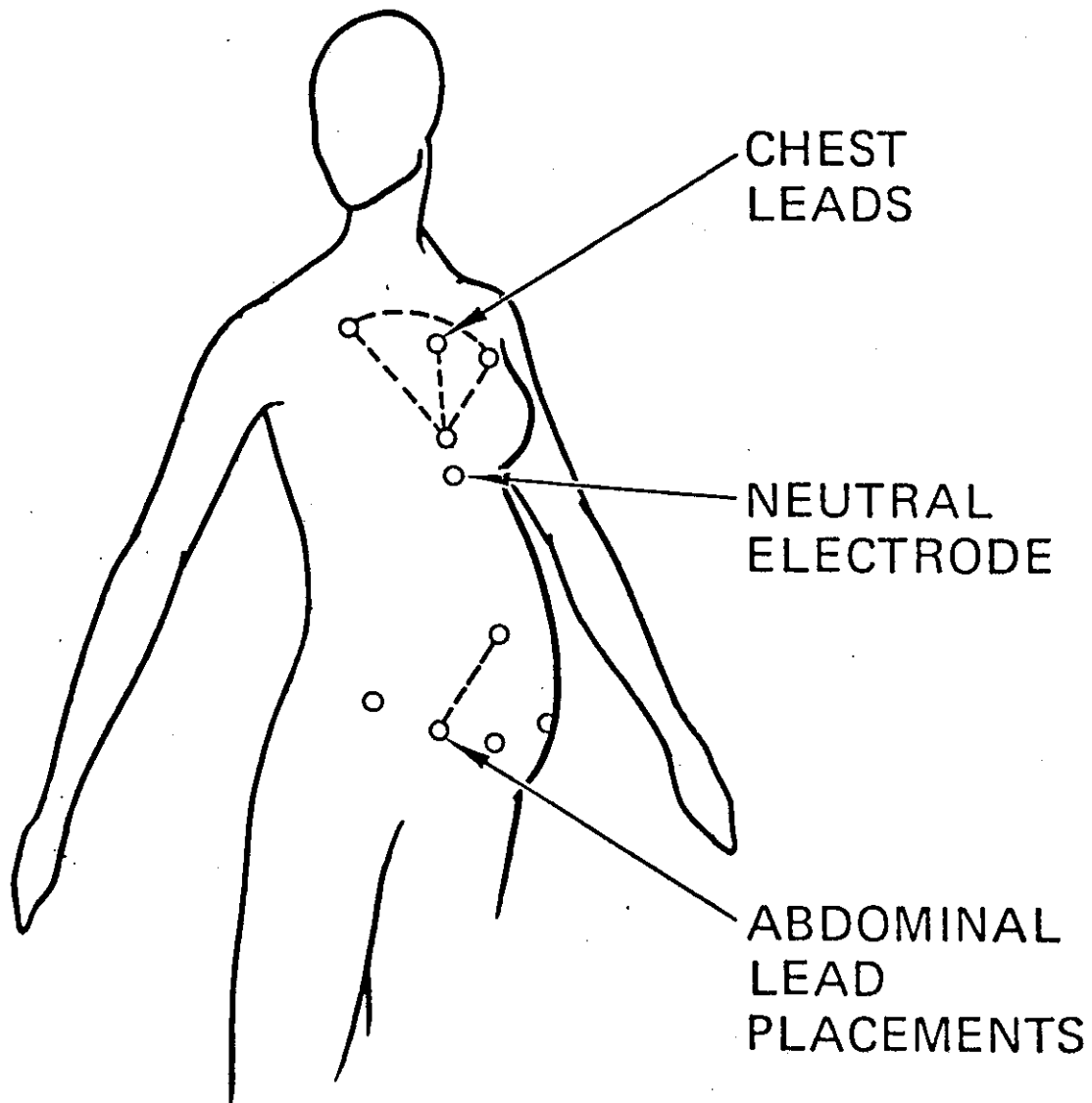
(c)



(a)



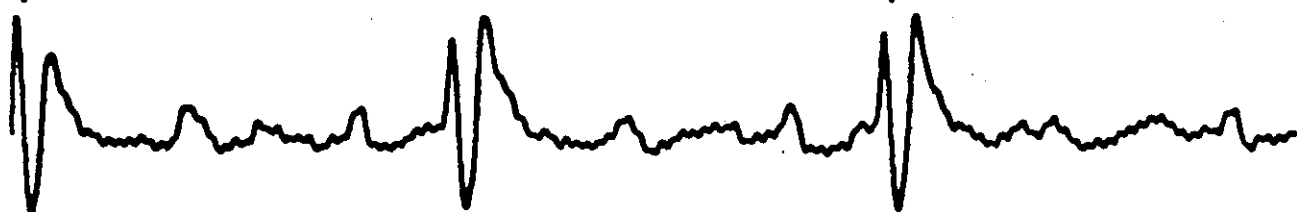
(b)



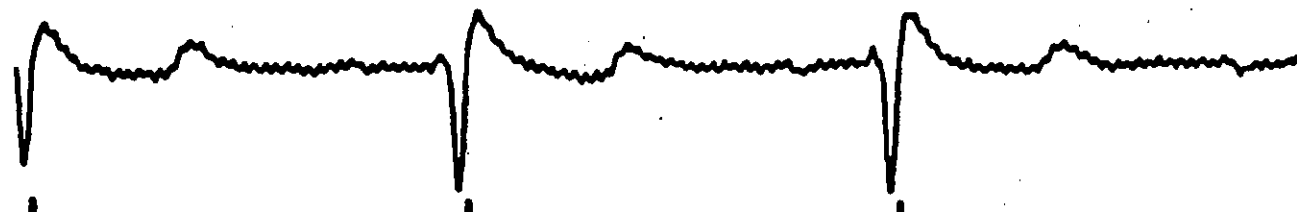




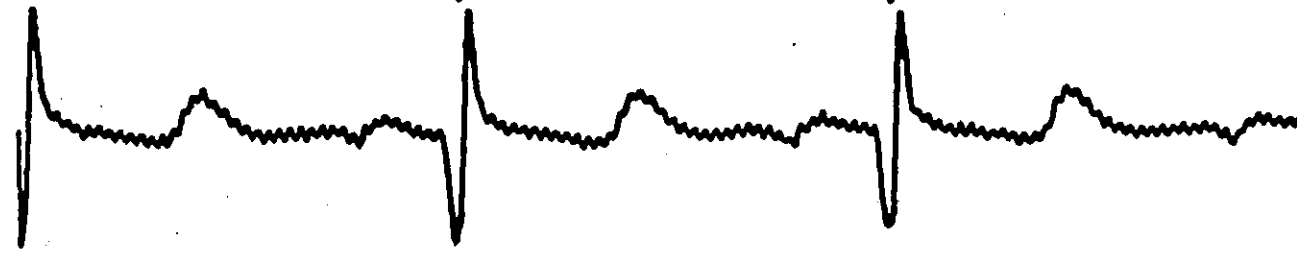
CHEST  
#1



#2



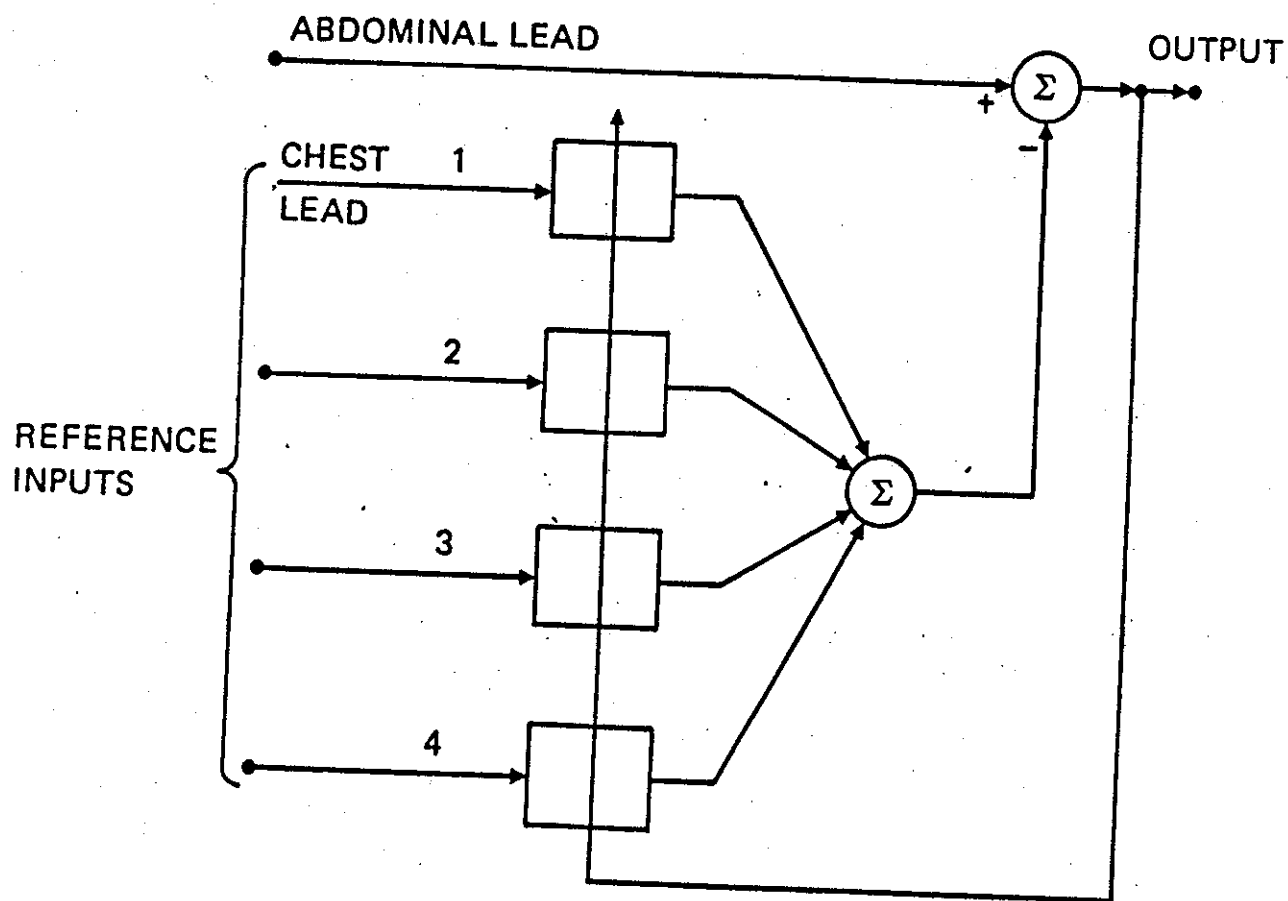
#3

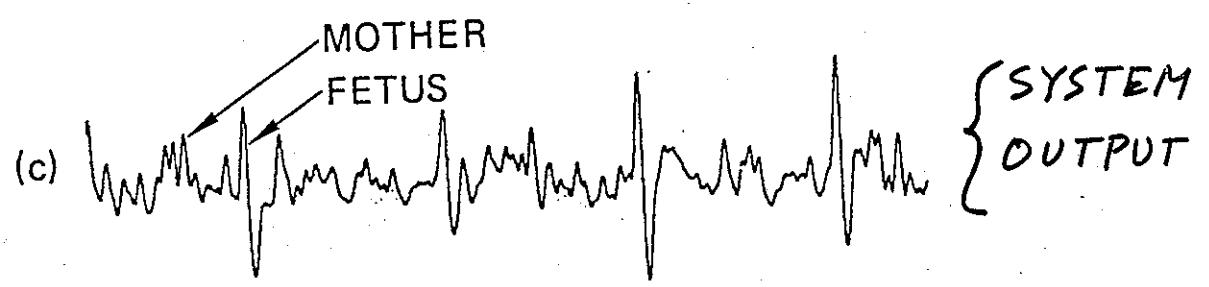
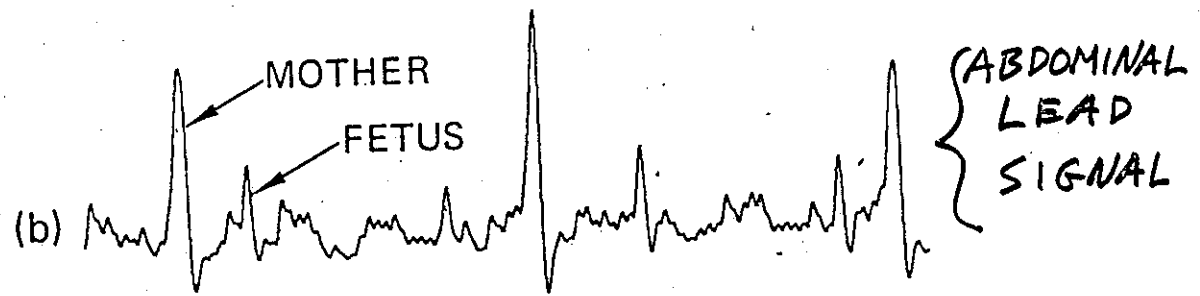
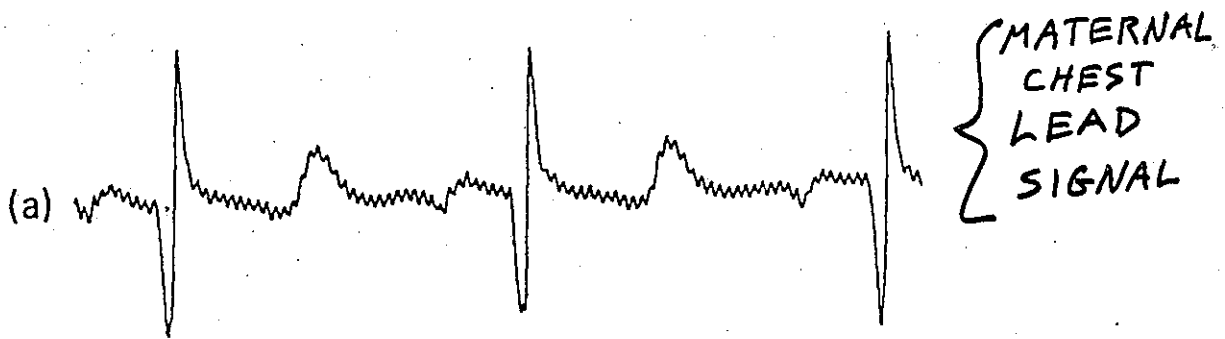


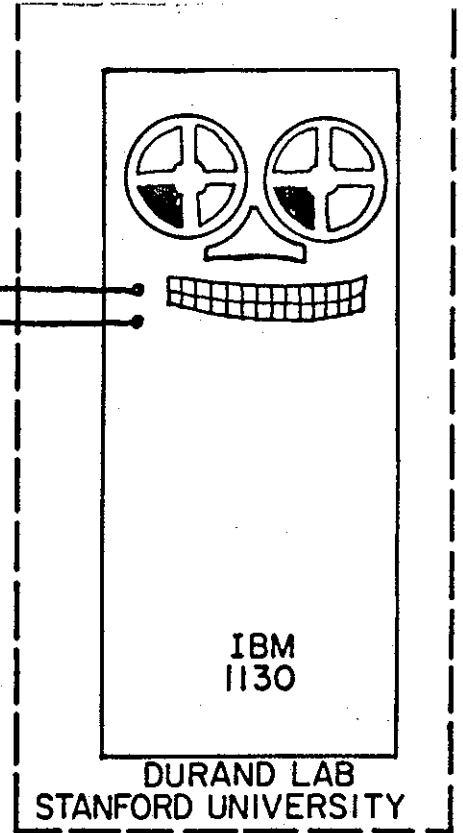
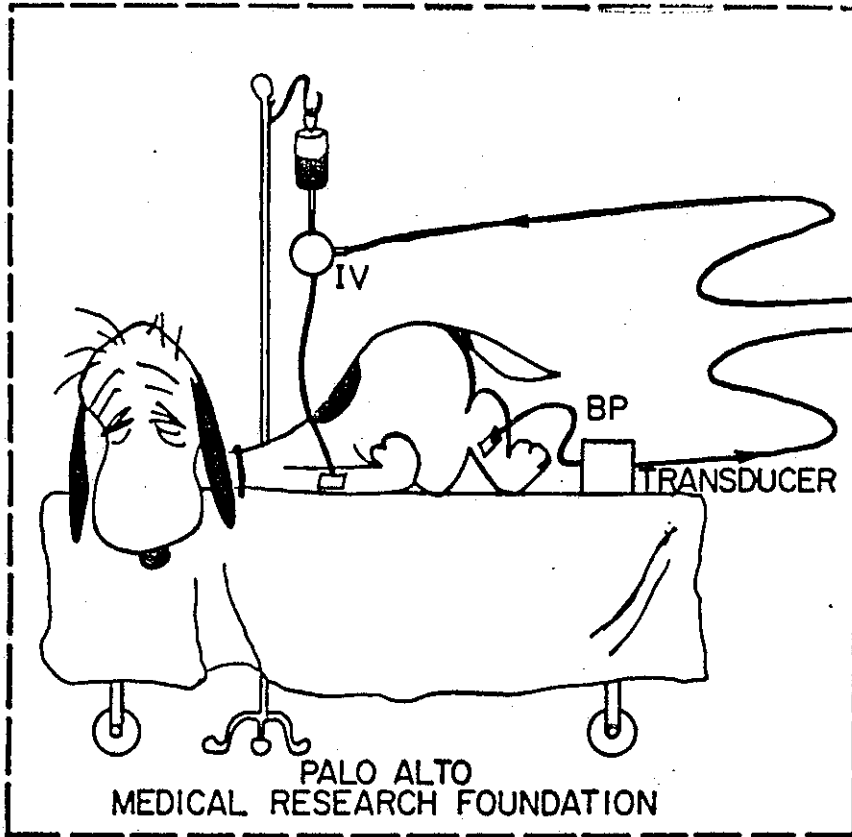
#4



BELLY







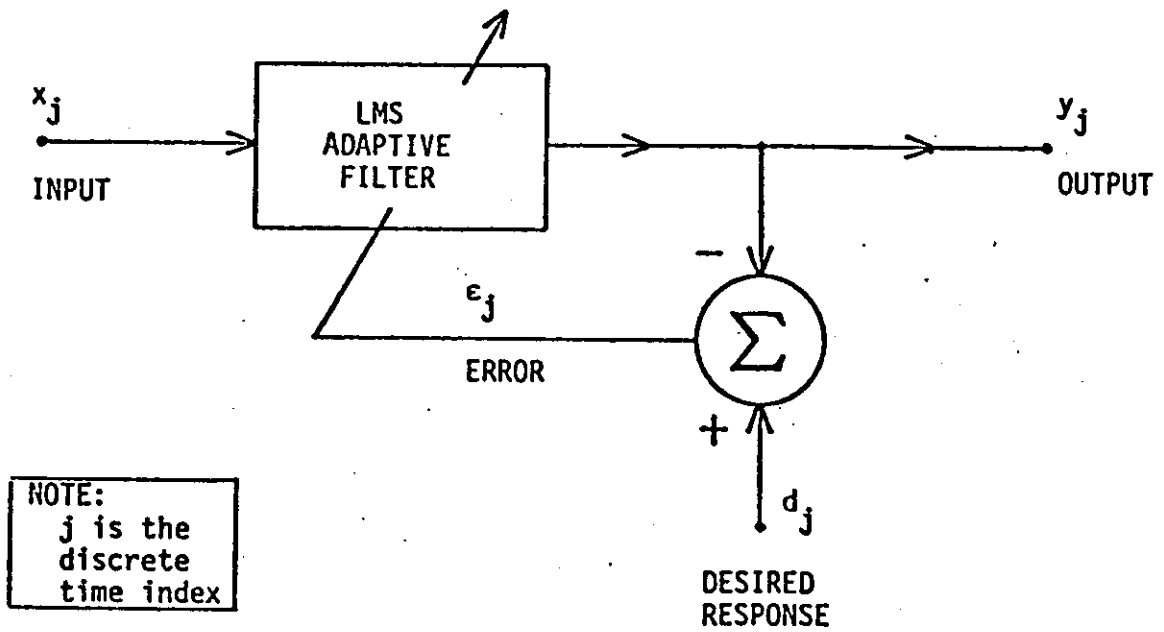


Fig. 1. Symbolic Representation of an Adaptive Transversal Filter Adapted by the LMS Algorithm.

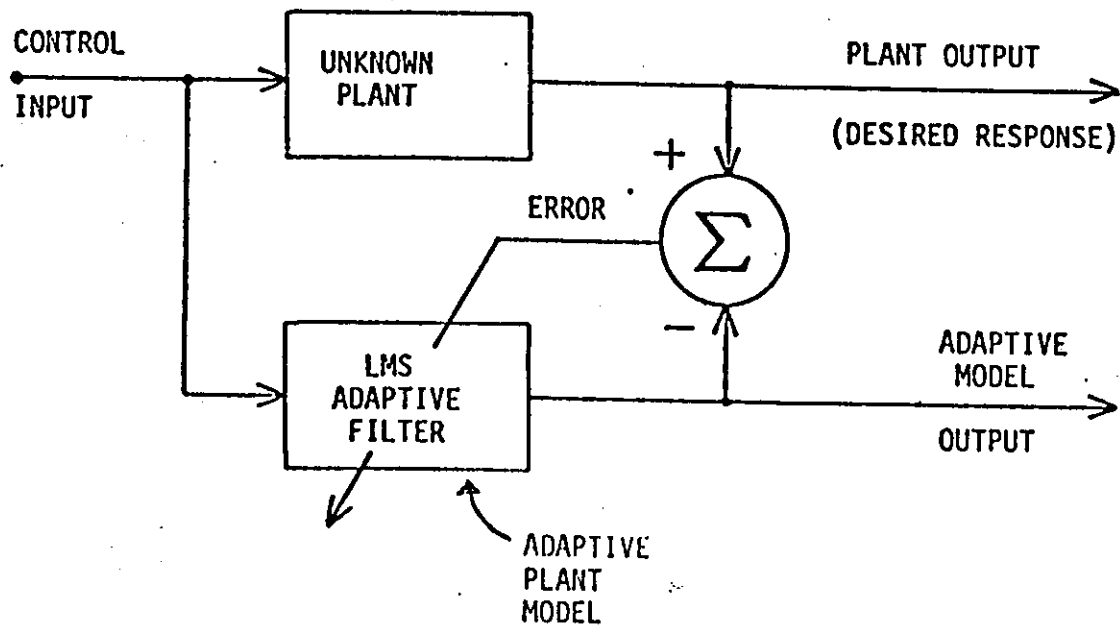


Fig. 2. Modeling an Unknown Plant by Means of an Adaptive Transversal Filter.

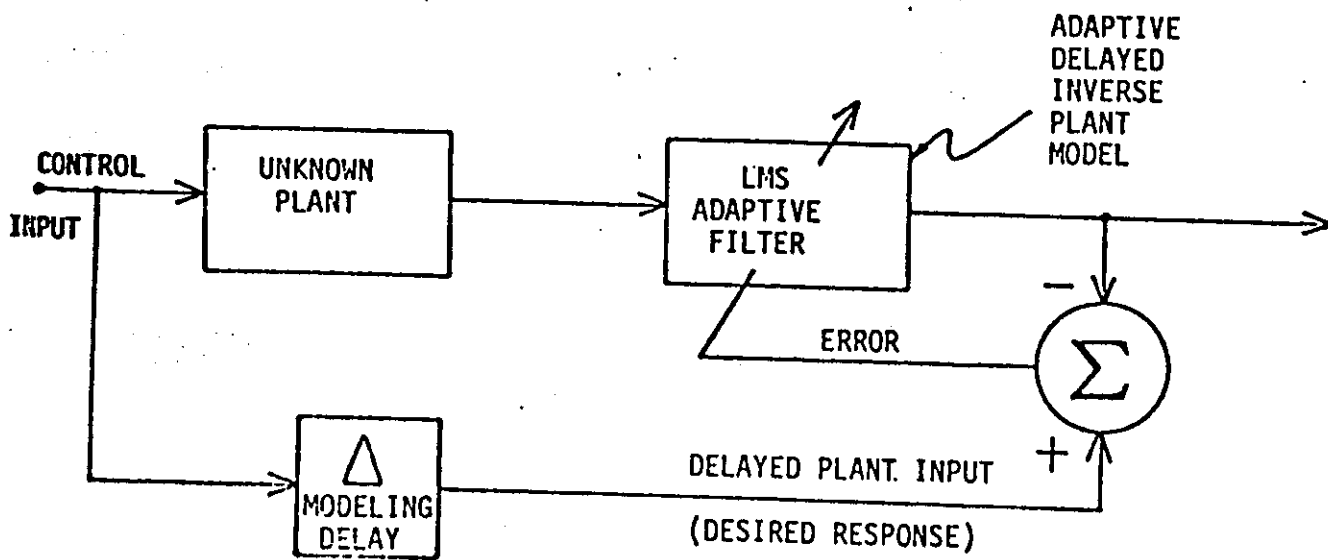


Fig. 4. Delayed Inverse Modeling of an Unknown Plant.

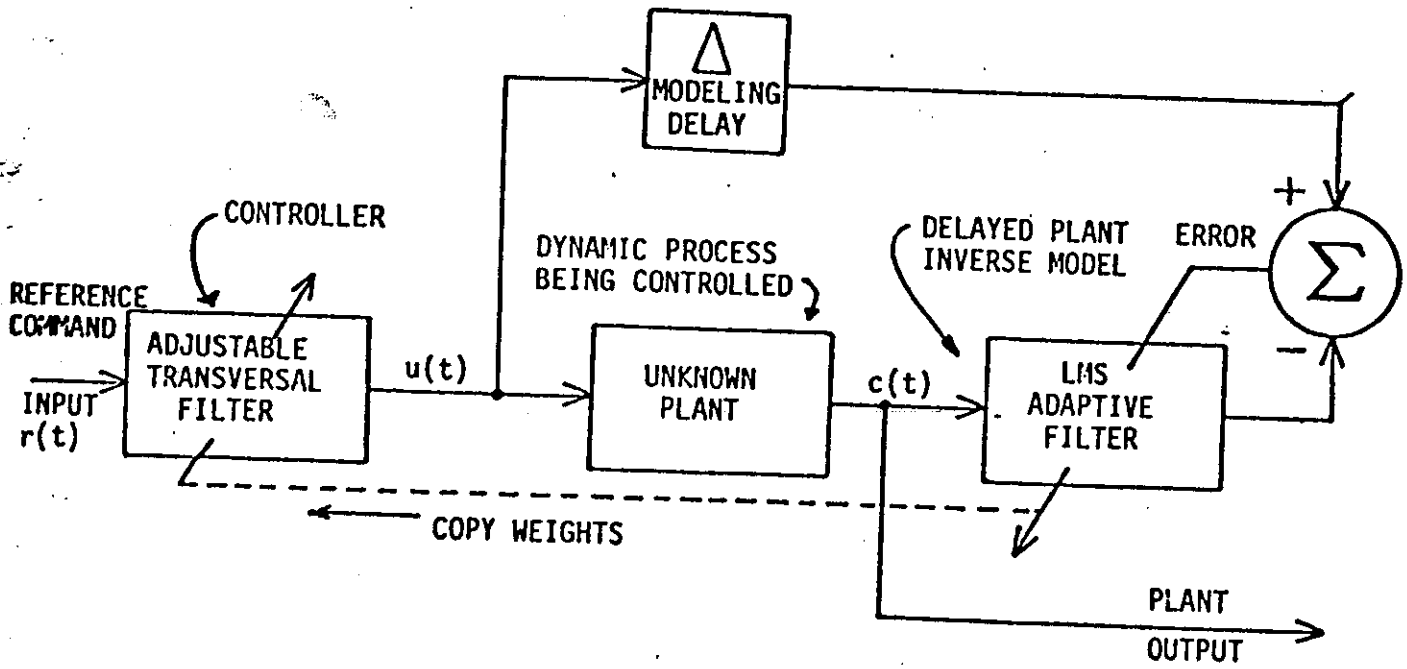


Fig. 5. Adaptive Inverse Model Control System.

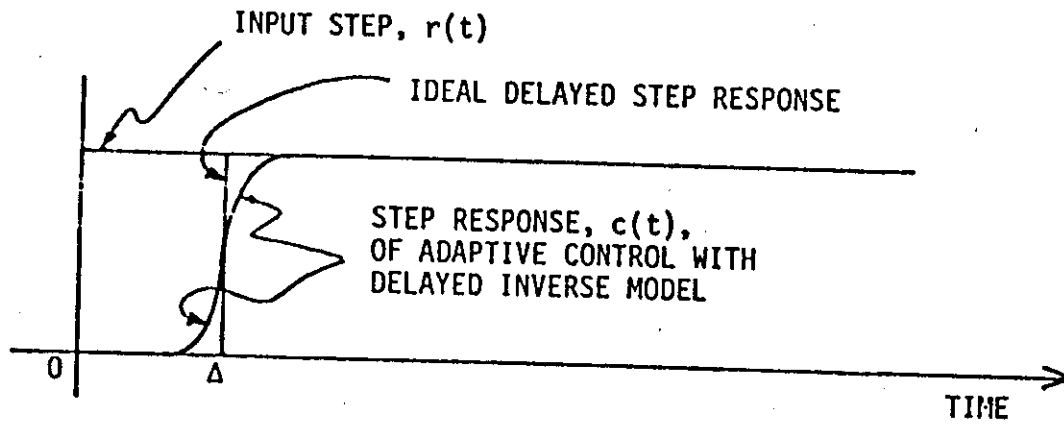


Fig. 6. Comparison of Ideal Step Response vs. Adaptive Delayed Inverse Model Control System Step Response.

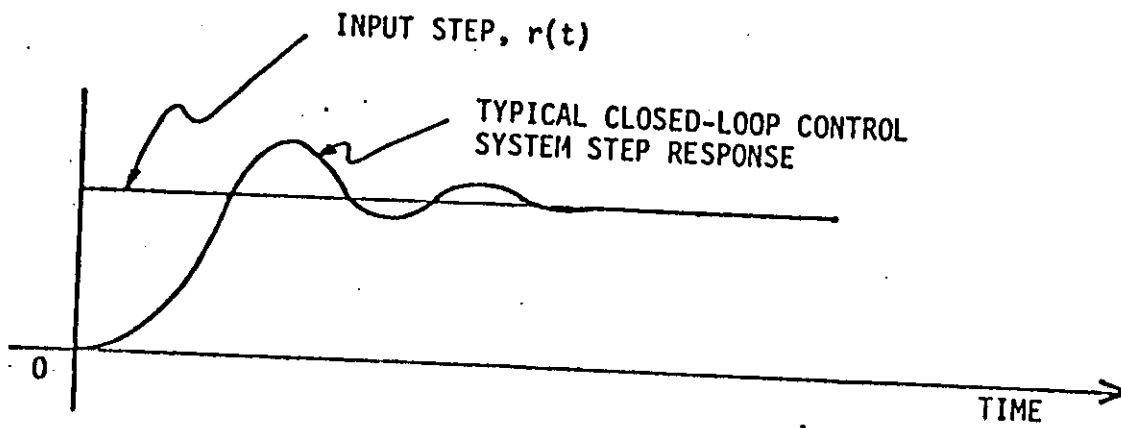


Fig. 7. Step Response of a Conventional Closed-Loop Control System.



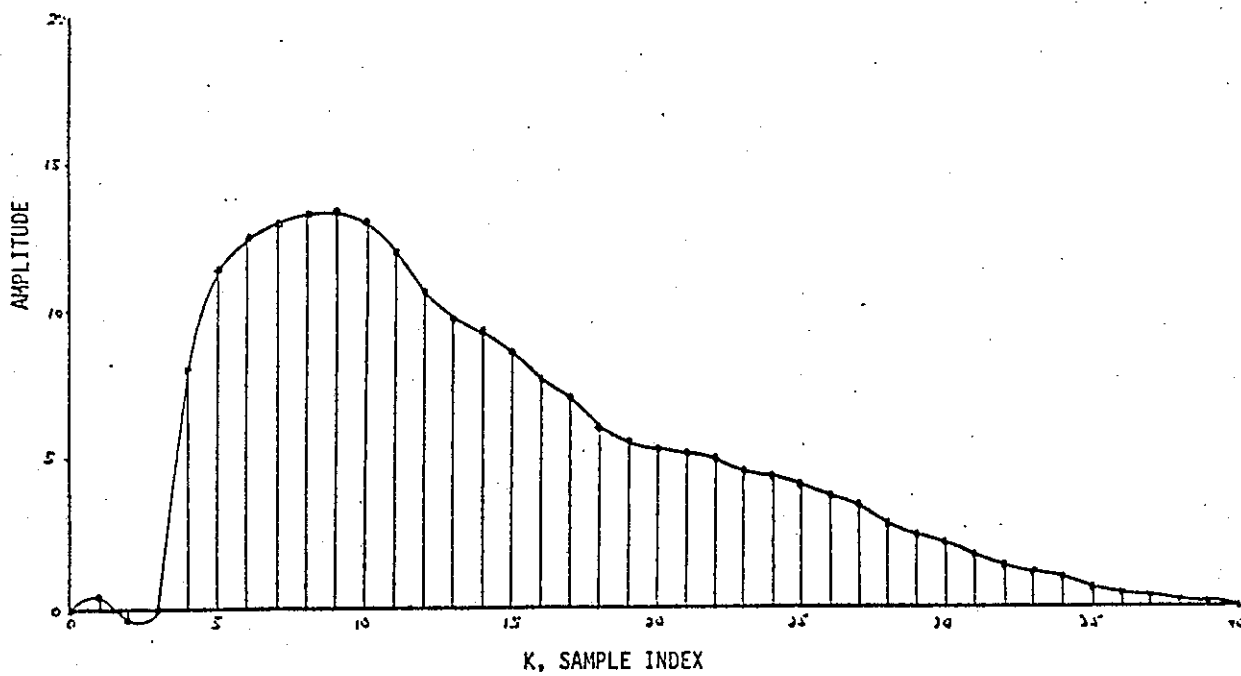


Fig. 2.5. A Discrete Impulse Response Whose Inverse Will Be Sought.

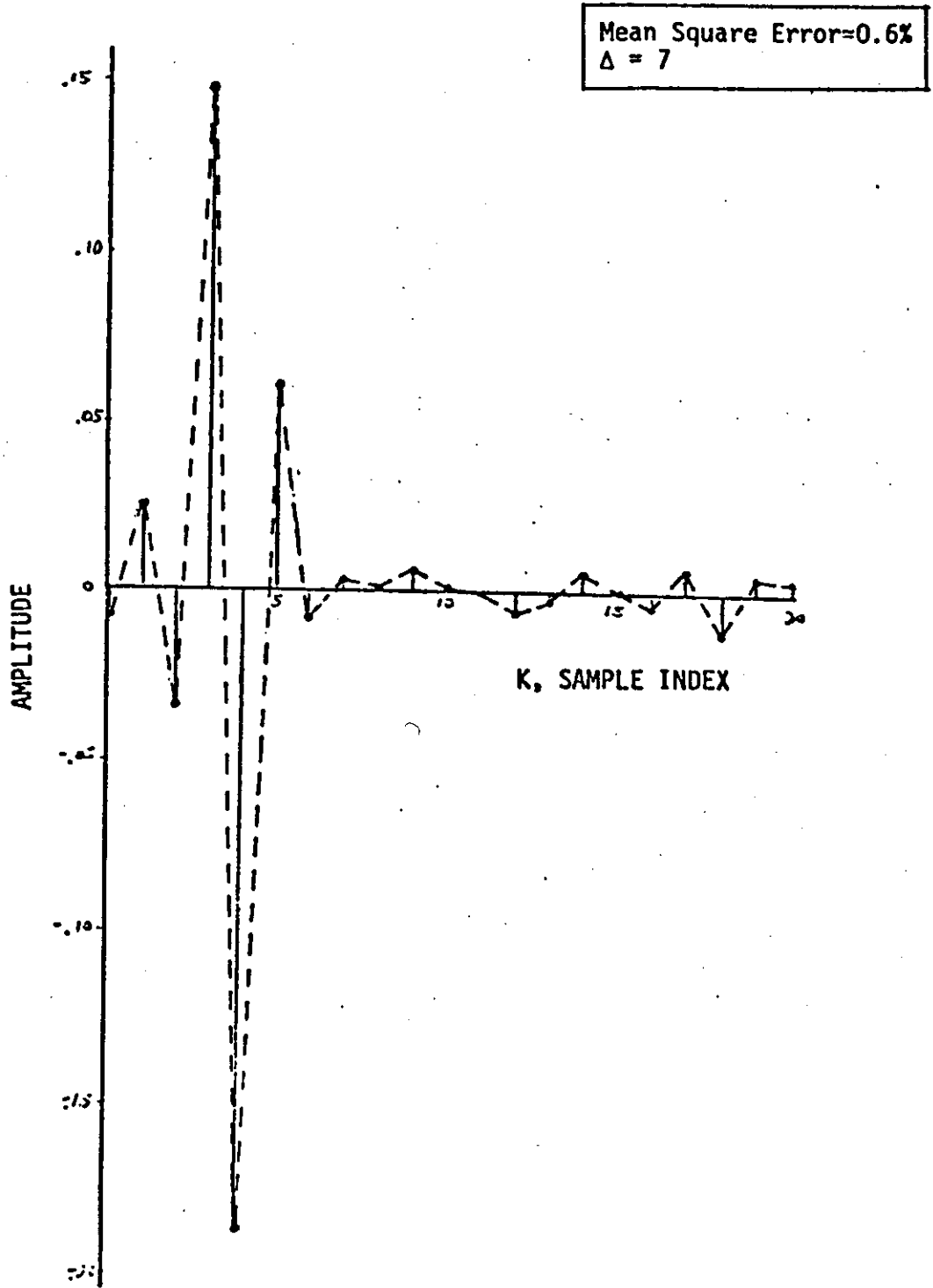


Fig. 2.7. A 21 Weight Delayed Approximate Inverse Impulse Response (7 Units of Delay).

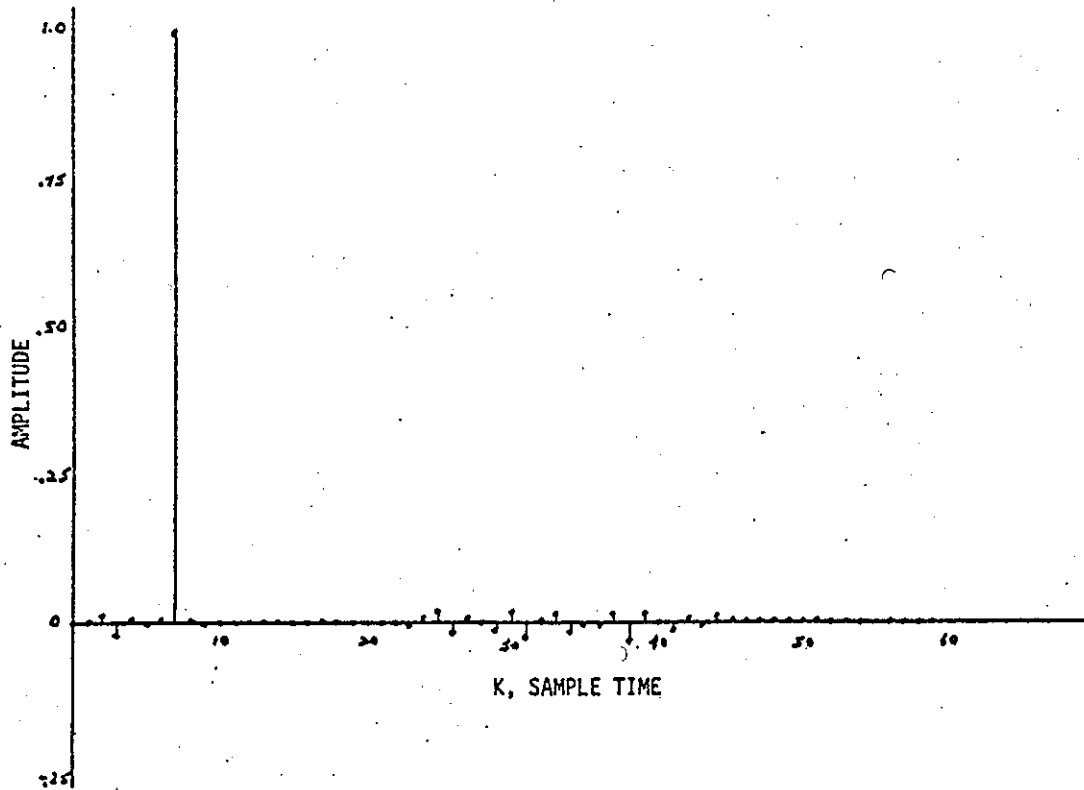


Fig. 2.8. Convolution of the Channel Impulse Response with Its Approximate Inverse

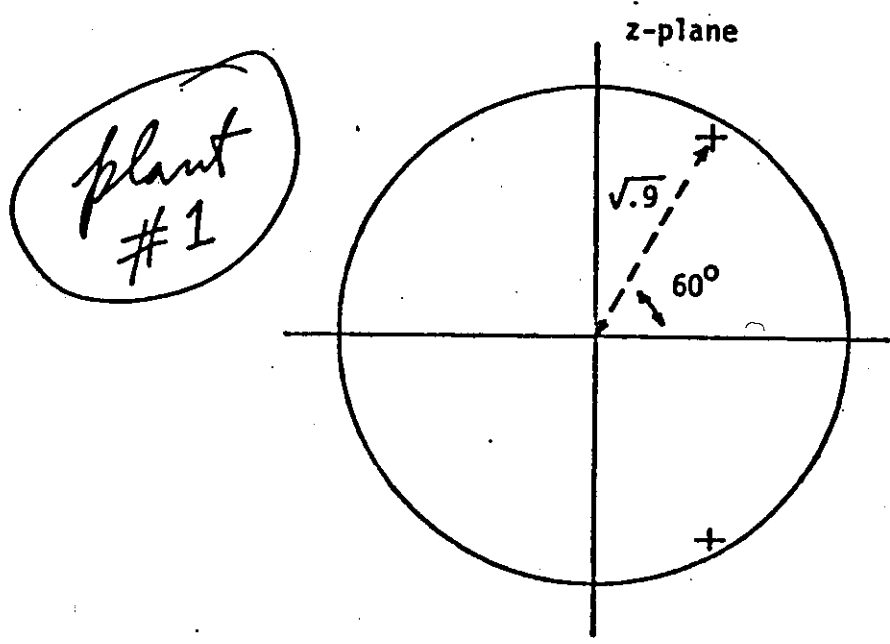


Fig. 10. Pole Locations of Plant #1.

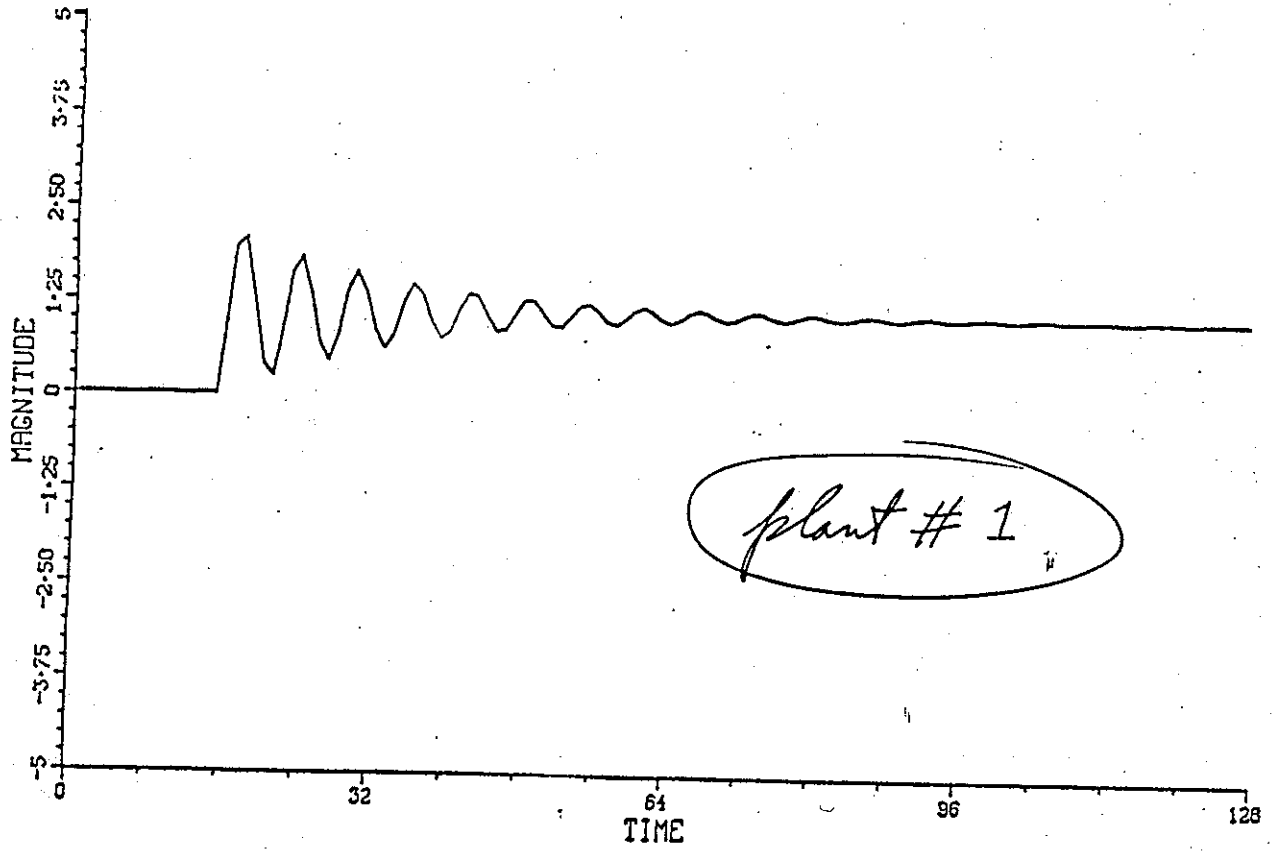


Fig. 11, Step Response of Plant #1.

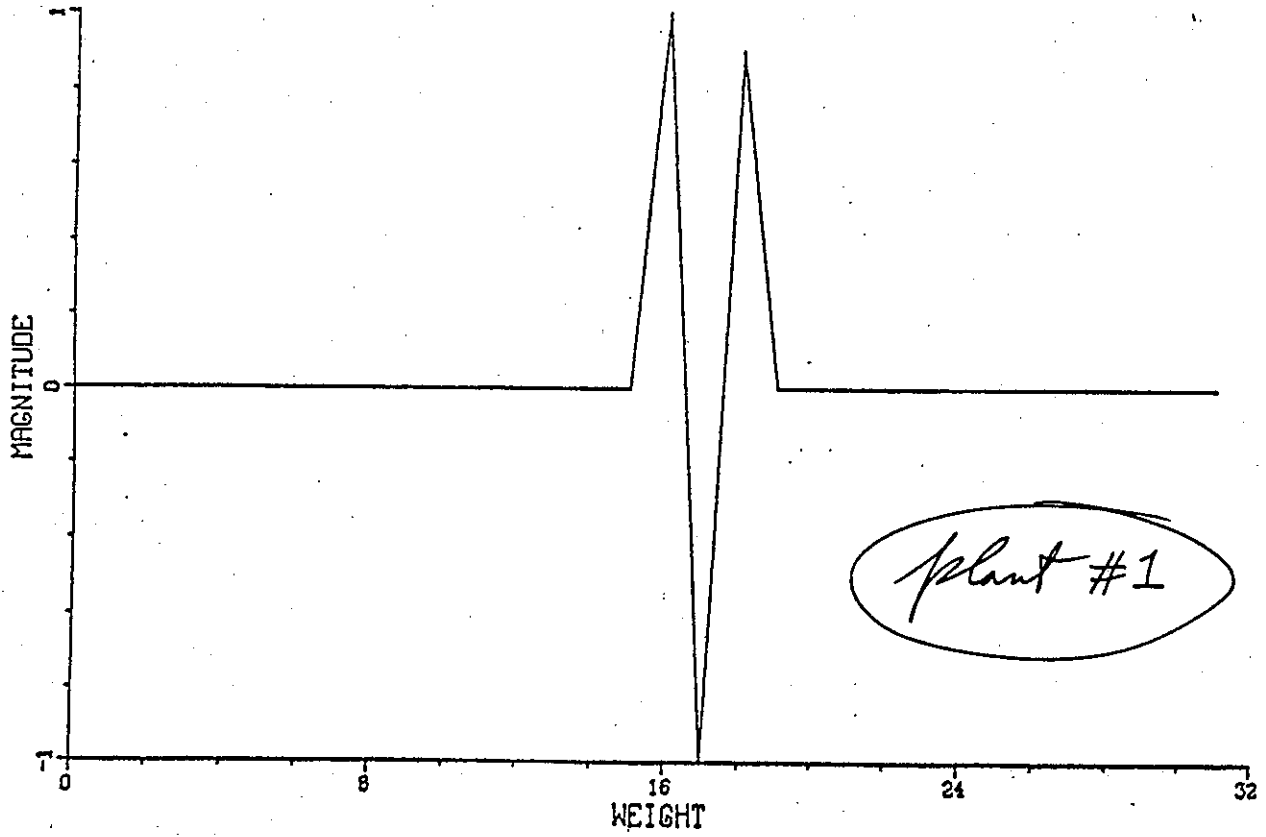


Fig. 12. Impulse Response of Inverse of Plant #1.

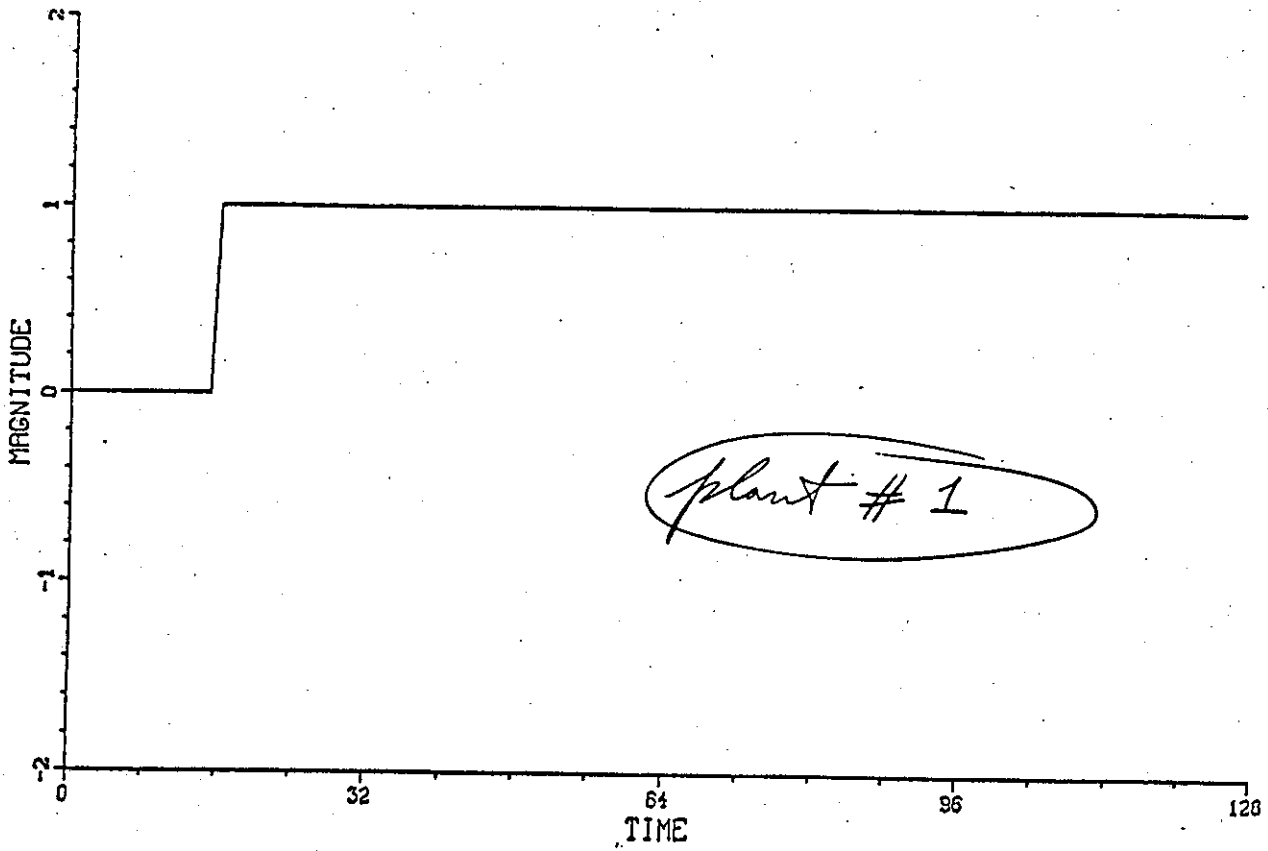


Fig. 13. Step Response of Inverse Model Control with Plant #1.

*Proc. IEEE, vol. 78, no. 9, pp 1415-1441, Sept., 1990*

# 30 Years of Adaptive Neural Networks: Perceptron, Madaline, and Backpropagation

Bernard Widrow

Michael A. Lehr

Stanford University Department of Electrical Engineering,

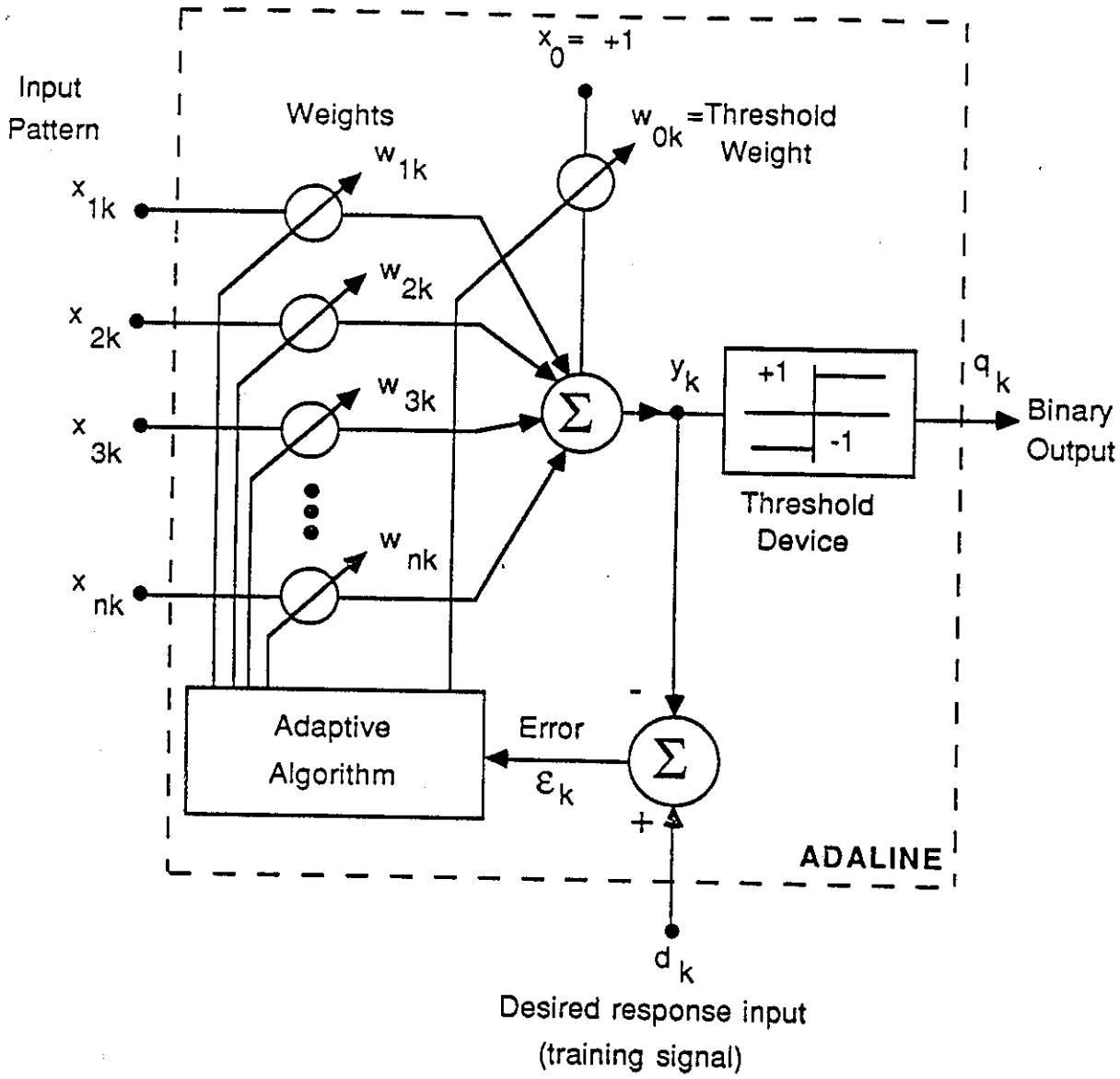
Stanford, CA 94305-4055

## Abstract

Fundamental developments in feedforward artificial neural networks from the past thirty years are reviewed. The central theme of this paper is a description of the history, origination, operating characteristics, and basic theory of several supervised neural network training algorithms including the Perceptron rule, the LMS algorithm, three Madaline rules, and the backpropagation technique. These methods were developed independently, but with the perspective of history they can all be related to each other. The concept which underlies these algorithms is the "minimal disturbance principle," which suggests that during training it is advisable to inject new information into a network in a manner which disturbs existing information to the smallest extent possible.

37a





$$X_k = [x_{0k}, x_{1k}, x_{2k}, \dots, x_{nk}]^T$$

$$= [+1, x_{1k}, x_{2k}, \dots, x_{nk}]^T$$

$$y_k = X_k^T W_k = W_k^T X_k$$

$$\epsilon_k = d_k - y_k$$

$$W_k = [w_{0k}, w_{1k}, w_{2k}, \dots, w_{nk}]^T$$

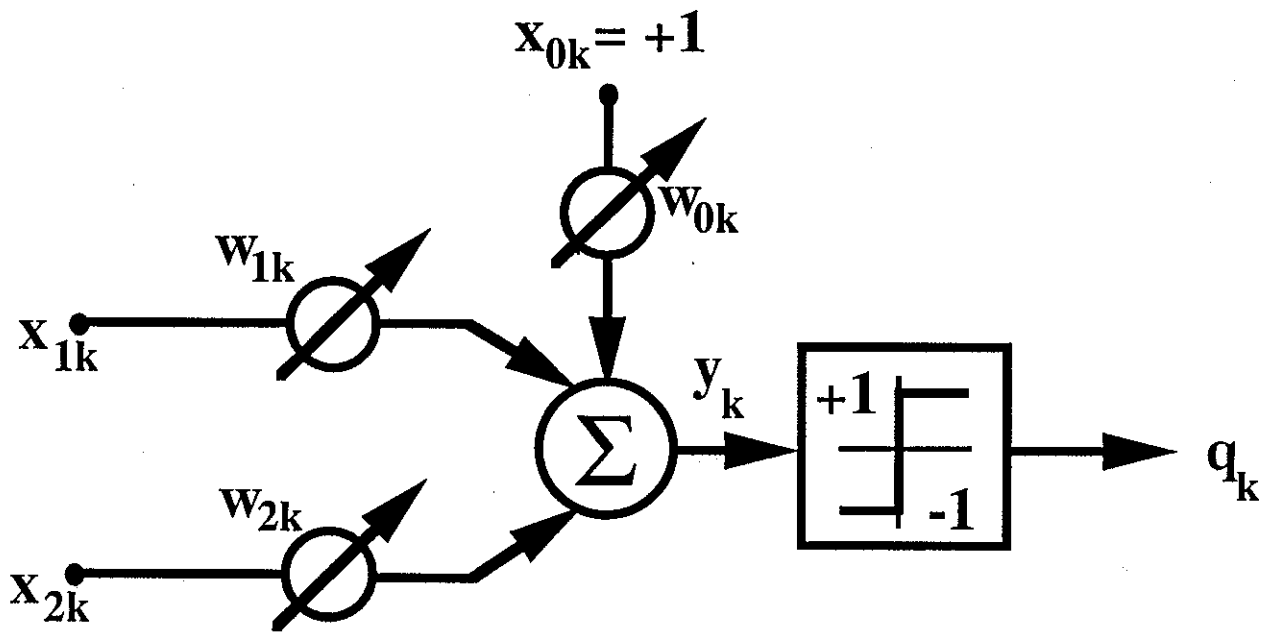
$k$  is time index

$T$  denotes vector transpose

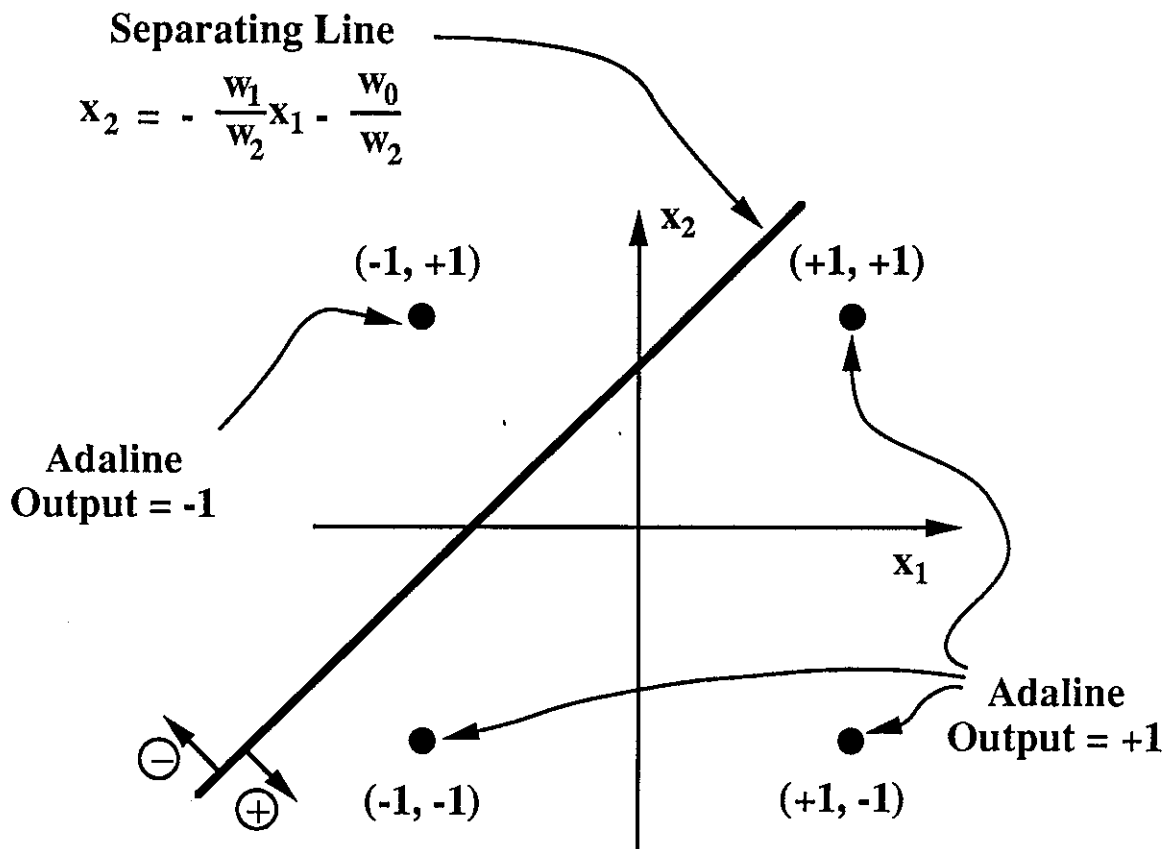
**An adaptive linear neuron (ADALINE).**

$$y = x_1 w_1 + x_2 w_2 + w_0 = 0$$

$$x_2 = -\frac{w_1}{w_2} x_1 - \frac{w_0}{w_2}$$



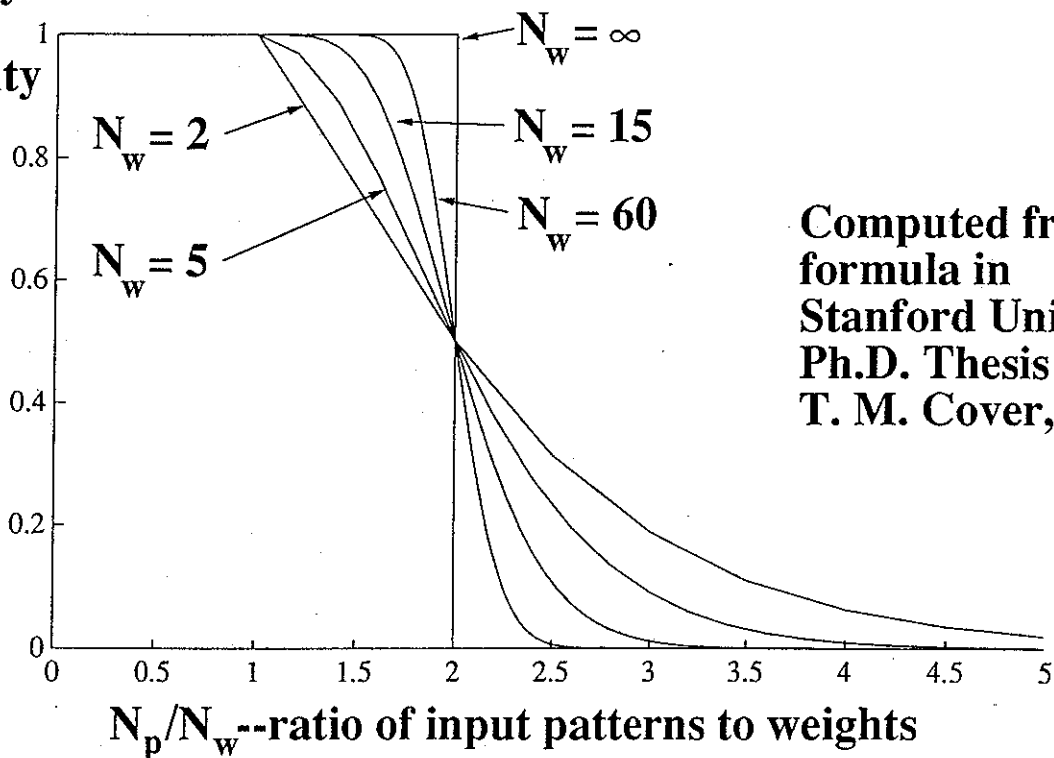
A two-Input Adaline



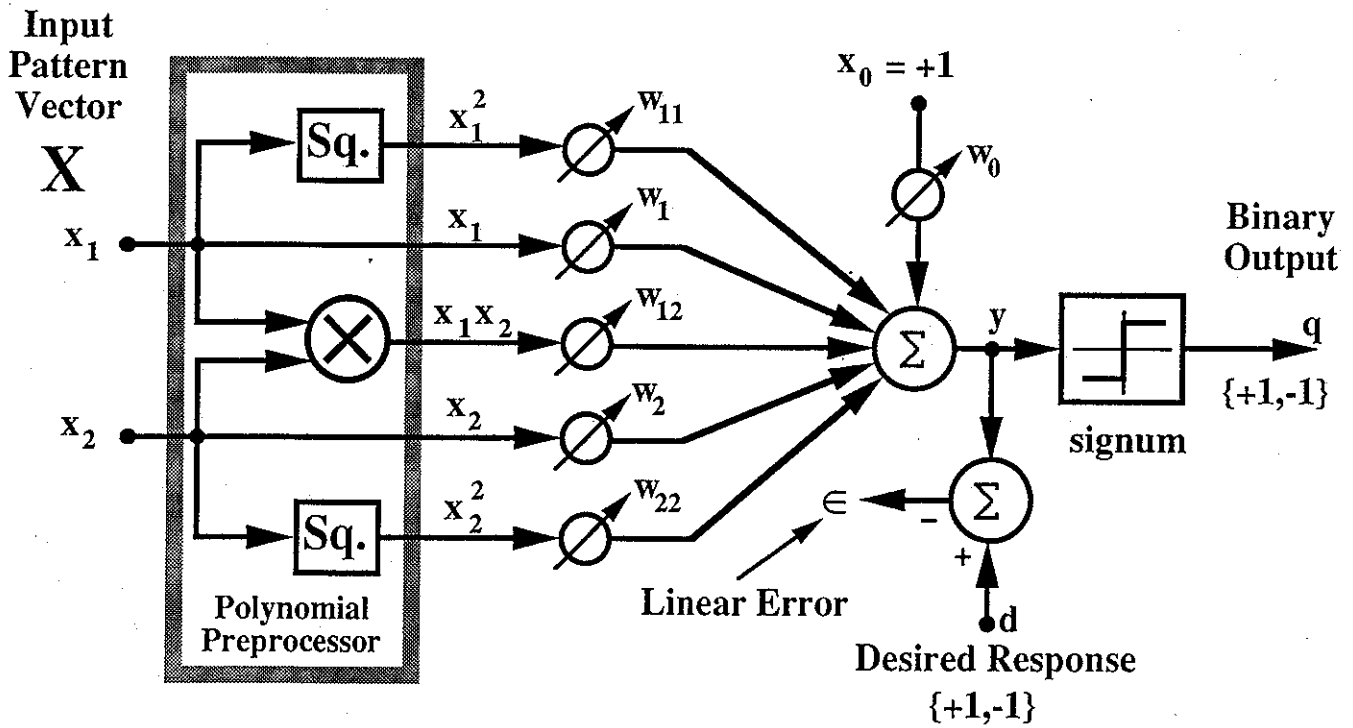
**Separating line in pattern space**

# Adaline Capacity

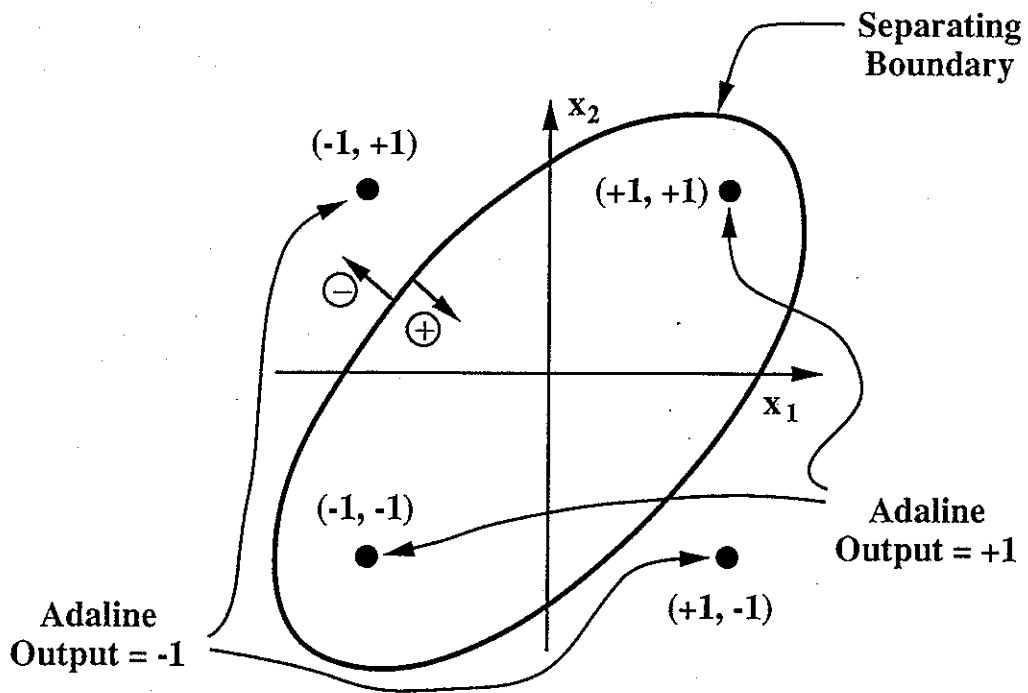
Probability  
of Linear  
Separability



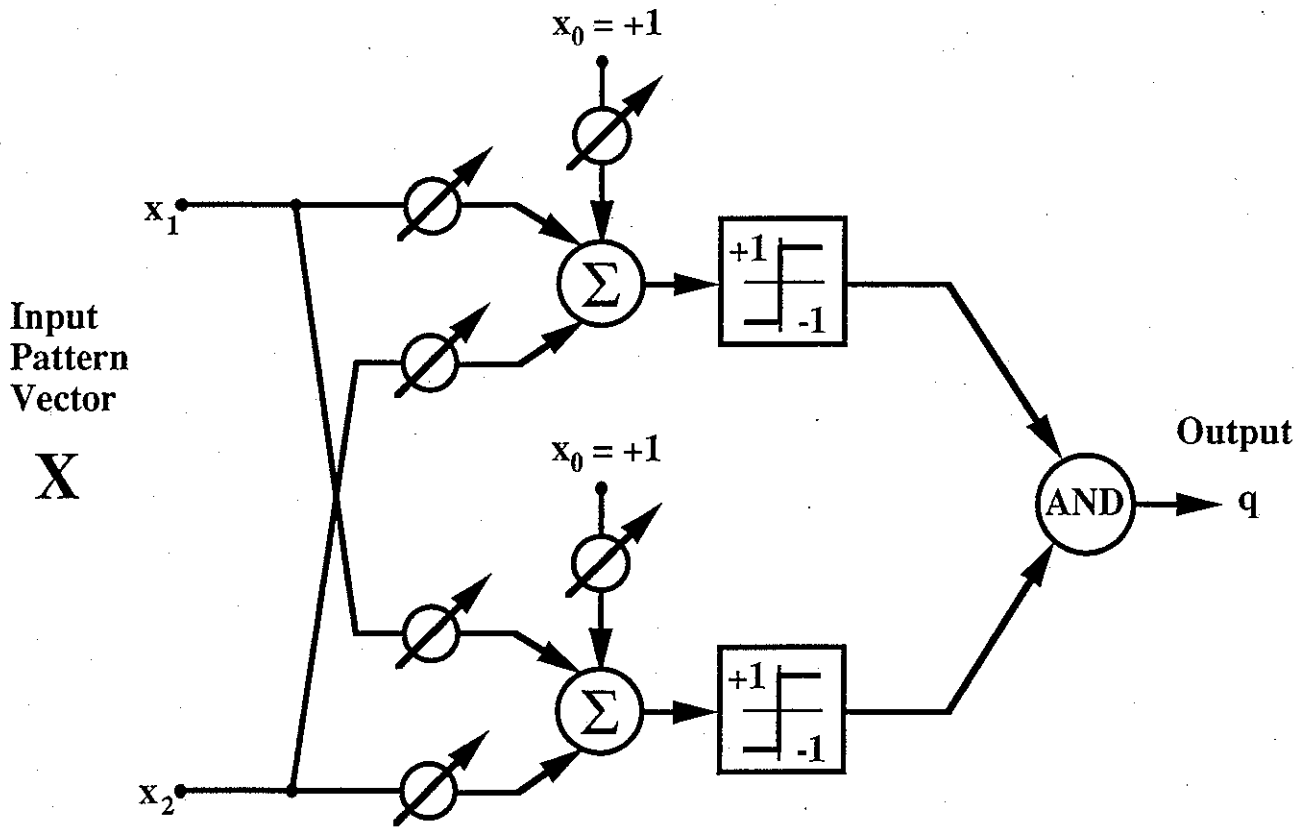
Computed from  
formula in  
Stanford Univ.  
Ph.D. Thesis by  
T. M. Cover, 1964



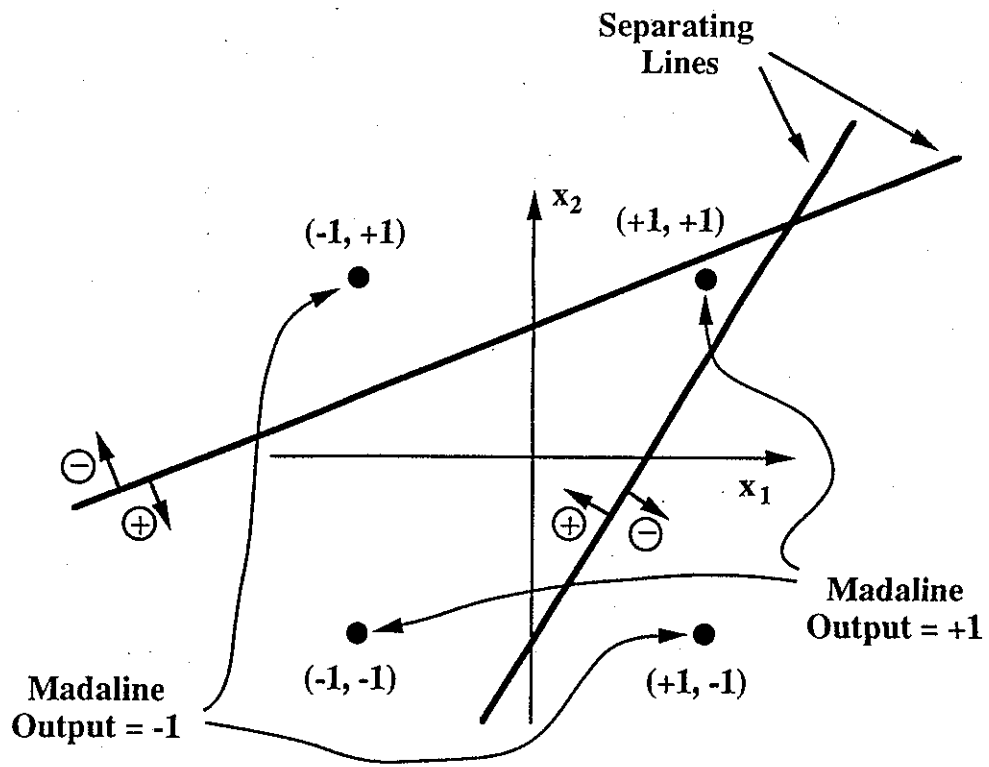
Adaline with polynomial preprocessor



**An elliptical separating boundary for the Exclusive NOR Function**

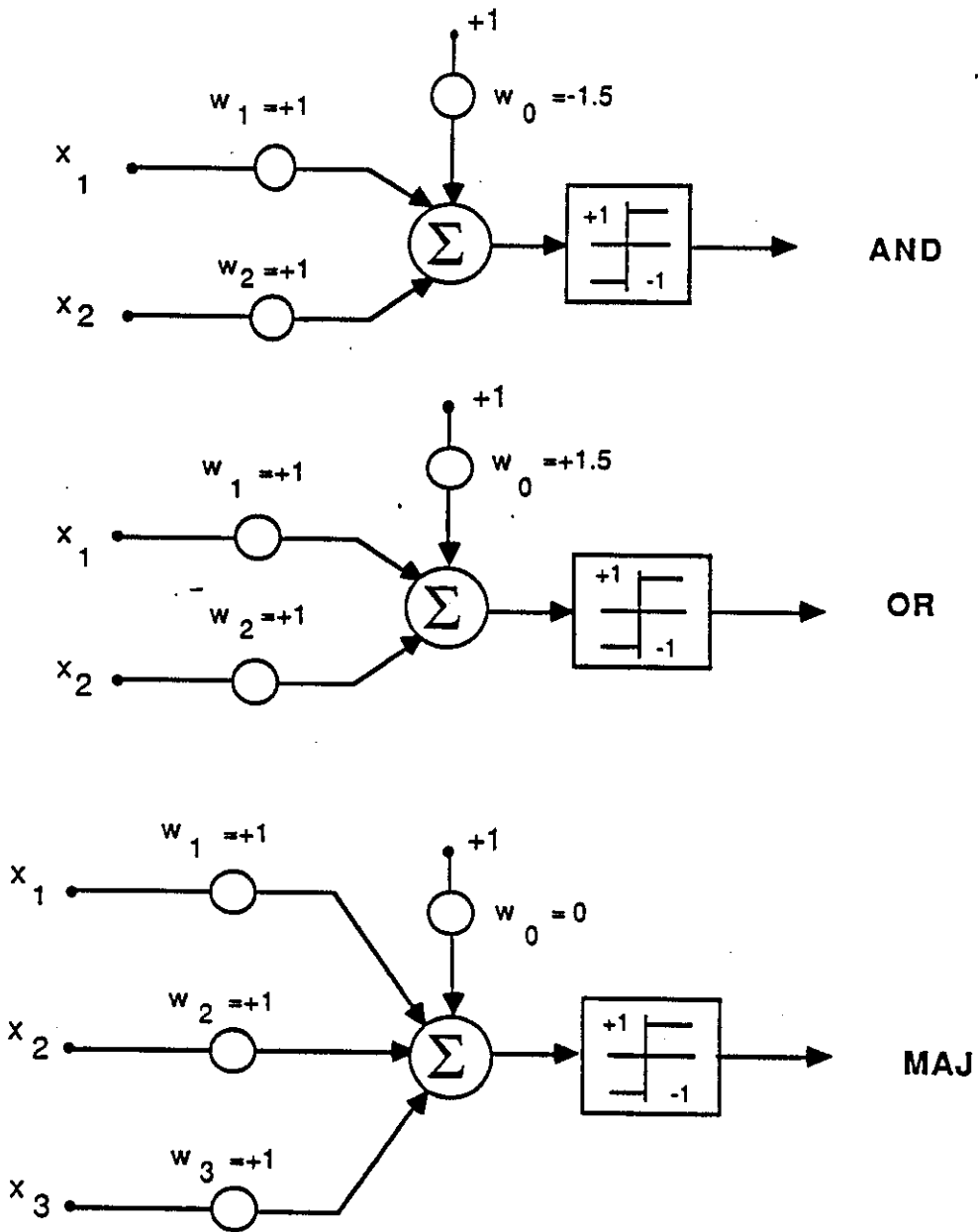


A two-Adaline form of Madaline I

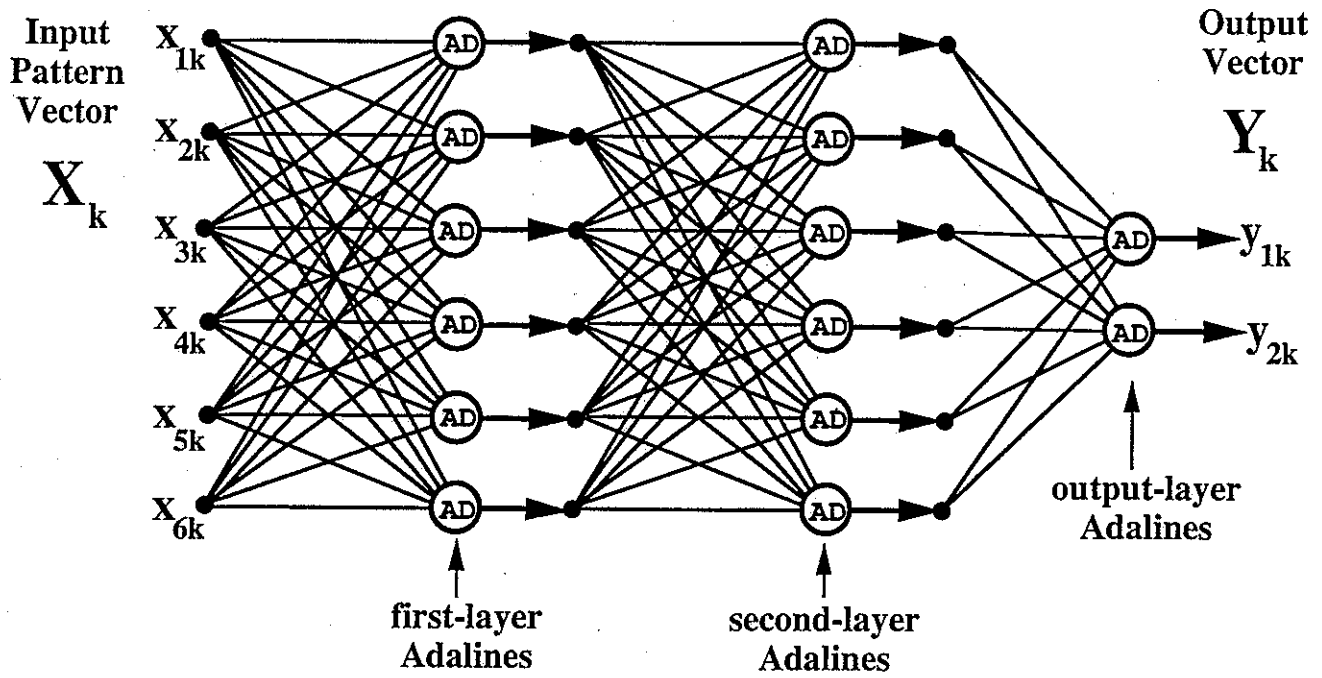


**Separating lines for the two-element Madaline**





A neuronal implementation of AND, OR, and MAJ logic functions.



**A three-layer adaptive neural network**

# Principle of Minimal Disturbance

*Adapt to reduce the output error for the current training pattern with minimal disturbance to the responses already learned.*

## $\alpha$ -LMS Algorithm

$$\epsilon_k \triangleq d_k - \mathbf{w}_k^T \mathbf{x}_k. \quad (1)$$

Changing the weights yields a corresponding change in the error:

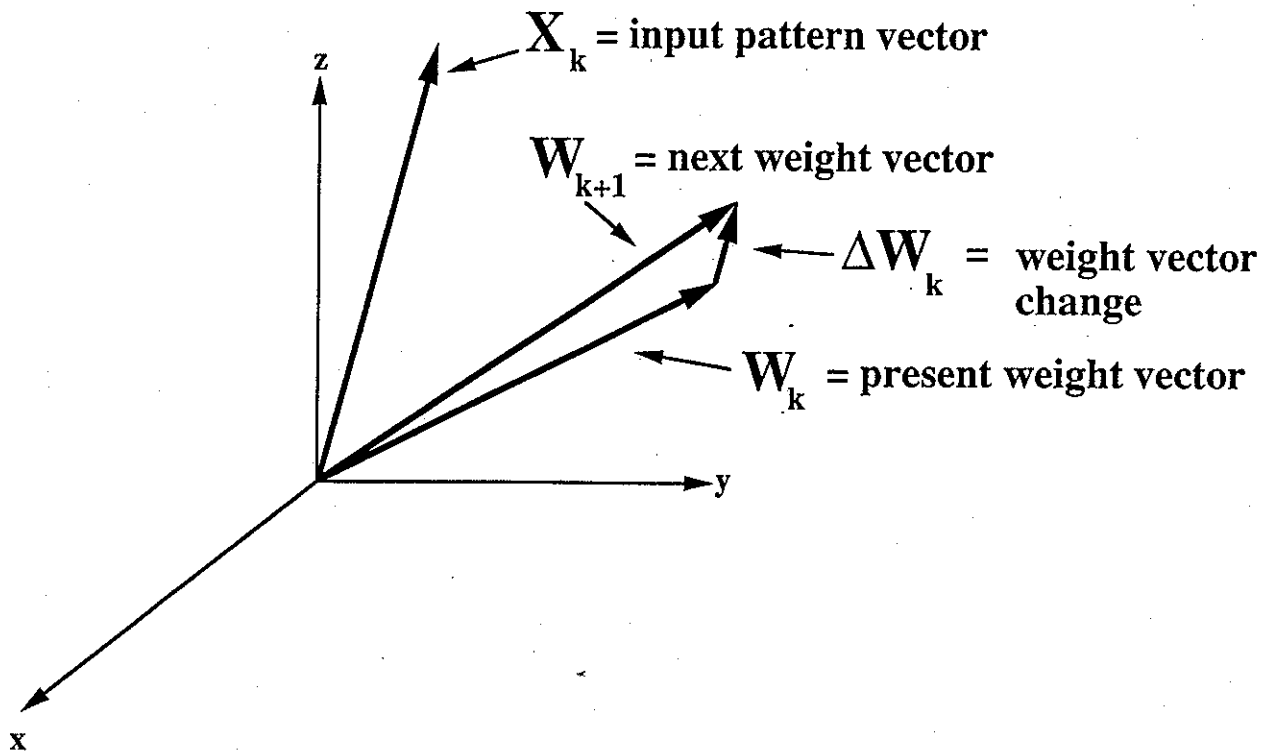
$$\Delta \epsilon_k = \Delta(d_k - \mathbf{w}_k^T \mathbf{x}_k) = -\mathbf{x}_k^T \Delta \mathbf{w}_k. \quad (2)$$

In accordance with the  $\alpha$ -LMS rule, the weight change is as follows:

$$\Delta \mathbf{w}_k = \mathbf{w}_{k+1} - \mathbf{w}_k = \alpha \frac{\epsilon_k \mathbf{x}_k}{|\mathbf{x}_k|^2}. \quad (3)$$

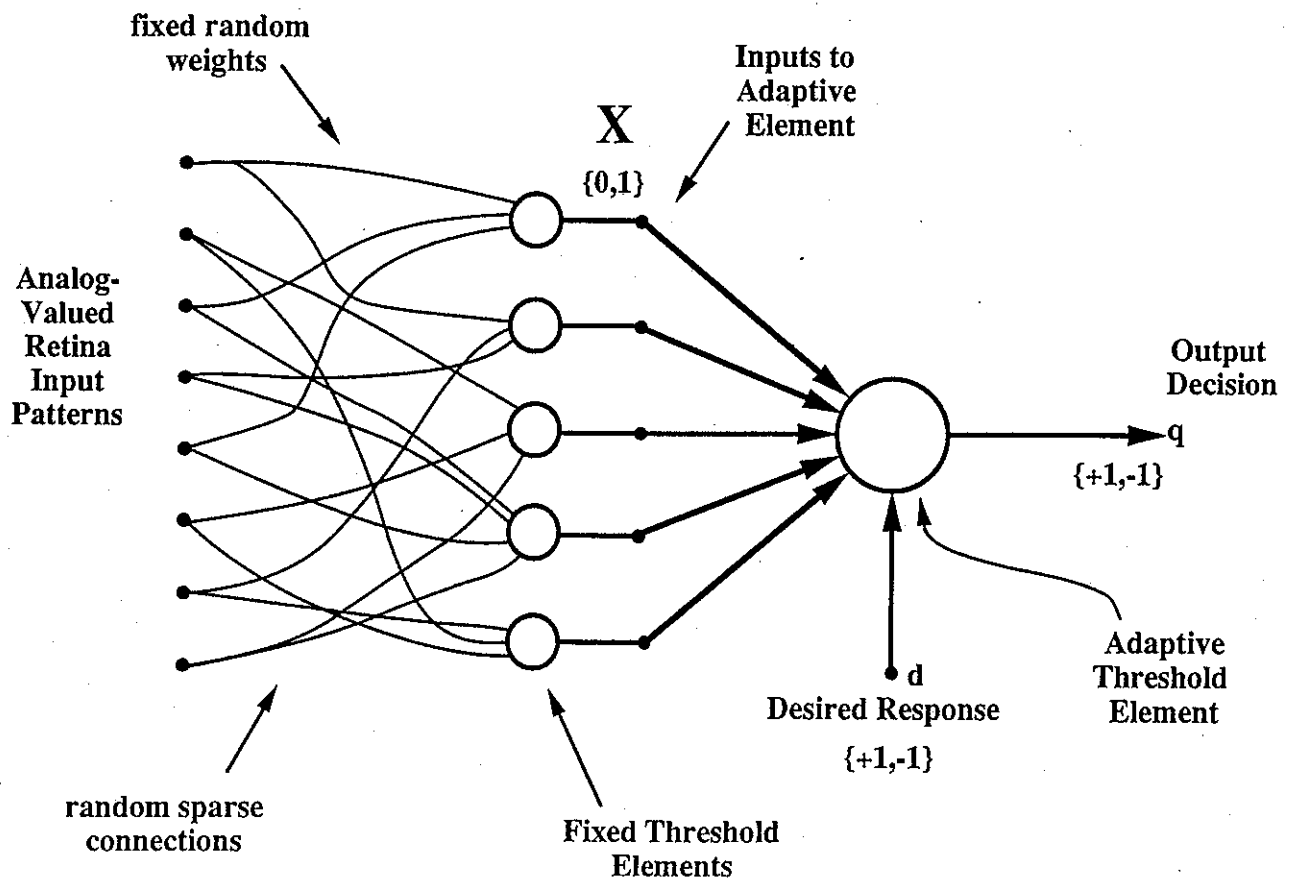
Combining Eqs. (2) and (3), we obtain

$$\Delta \epsilon_k = -\alpha \frac{\epsilon_k \mathbf{x}_k^T \mathbf{x}_k}{|\mathbf{x}_k|^2} = -\alpha \epsilon_k. \quad (4)$$

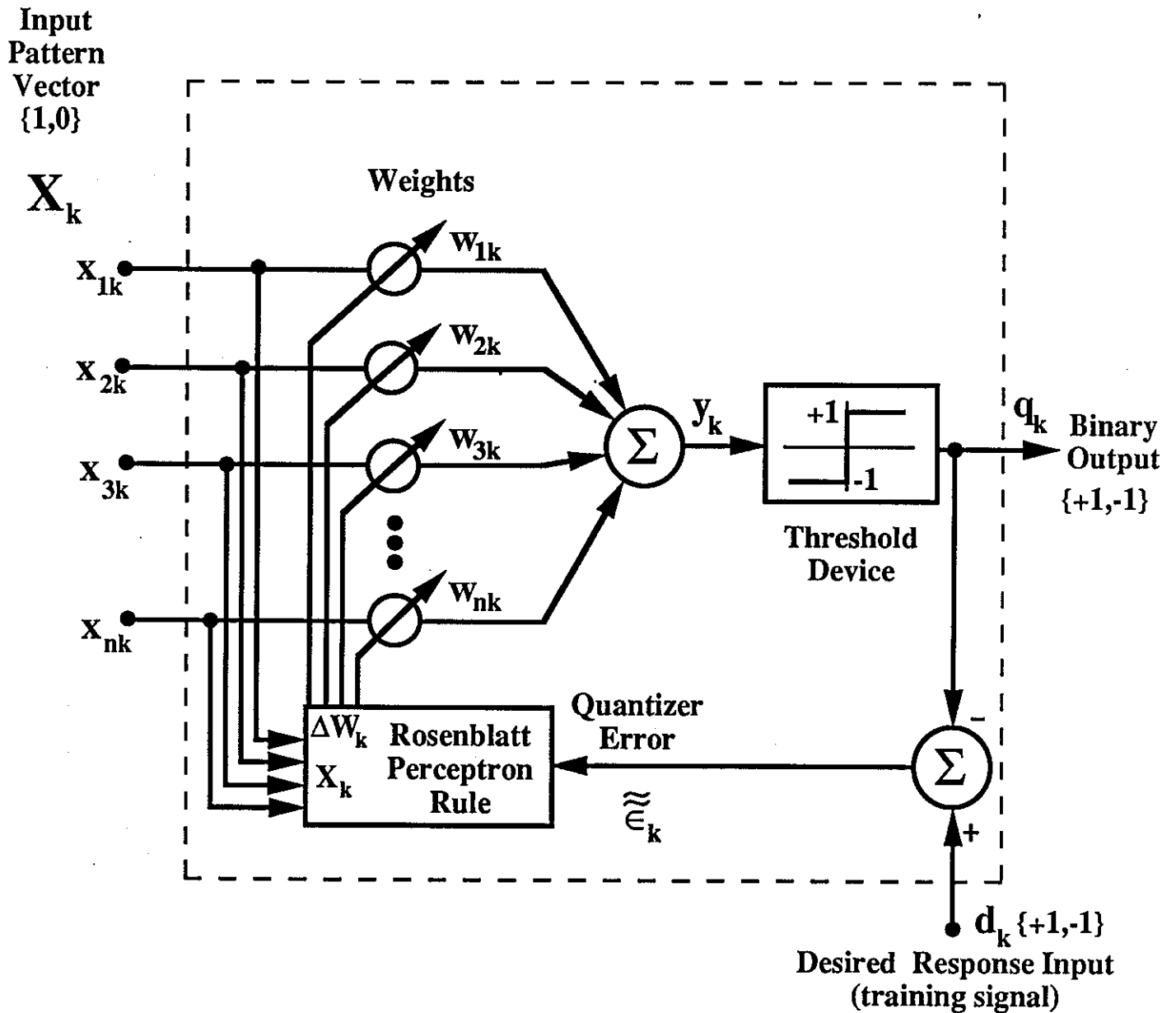


**Weight correction by the LMS rule**

# Rosenblatt's Perceptron



# Adaptive Threshold Element in the Perceptron



**Perceptron Rule:**

$$W_{k+1} = W_k + \mu \tilde{\epsilon}_k X_k$$

$\mu$  normally set to 1/2

# Perceptron Rule

- If response is OK, do not adapt weights.
- Otherwise adapt weights by a fixed distance along the  $X$ -Vector to reduce error

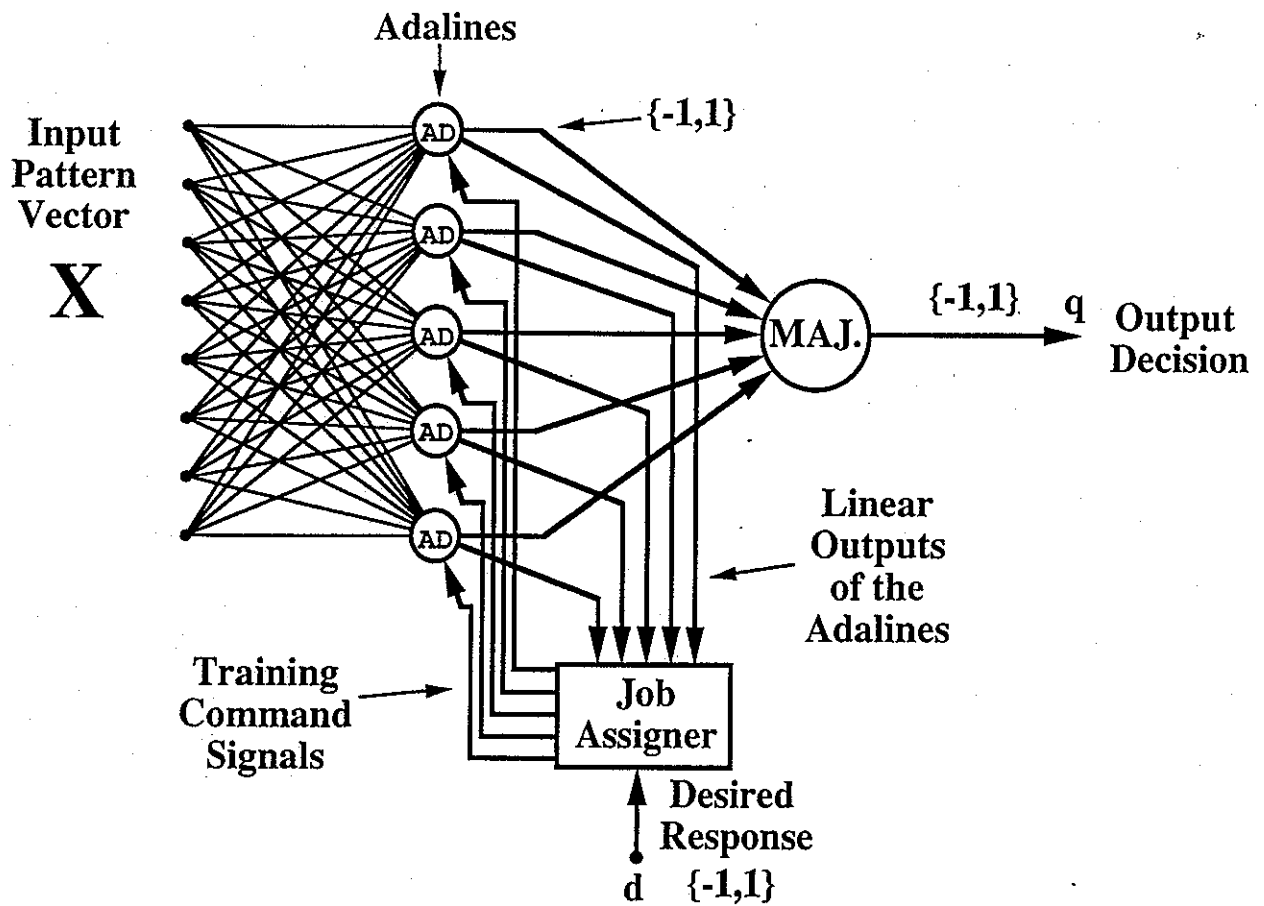
## Good Features

- Guaranteed to converge to solution if problem is linearly separable

## Bad Features

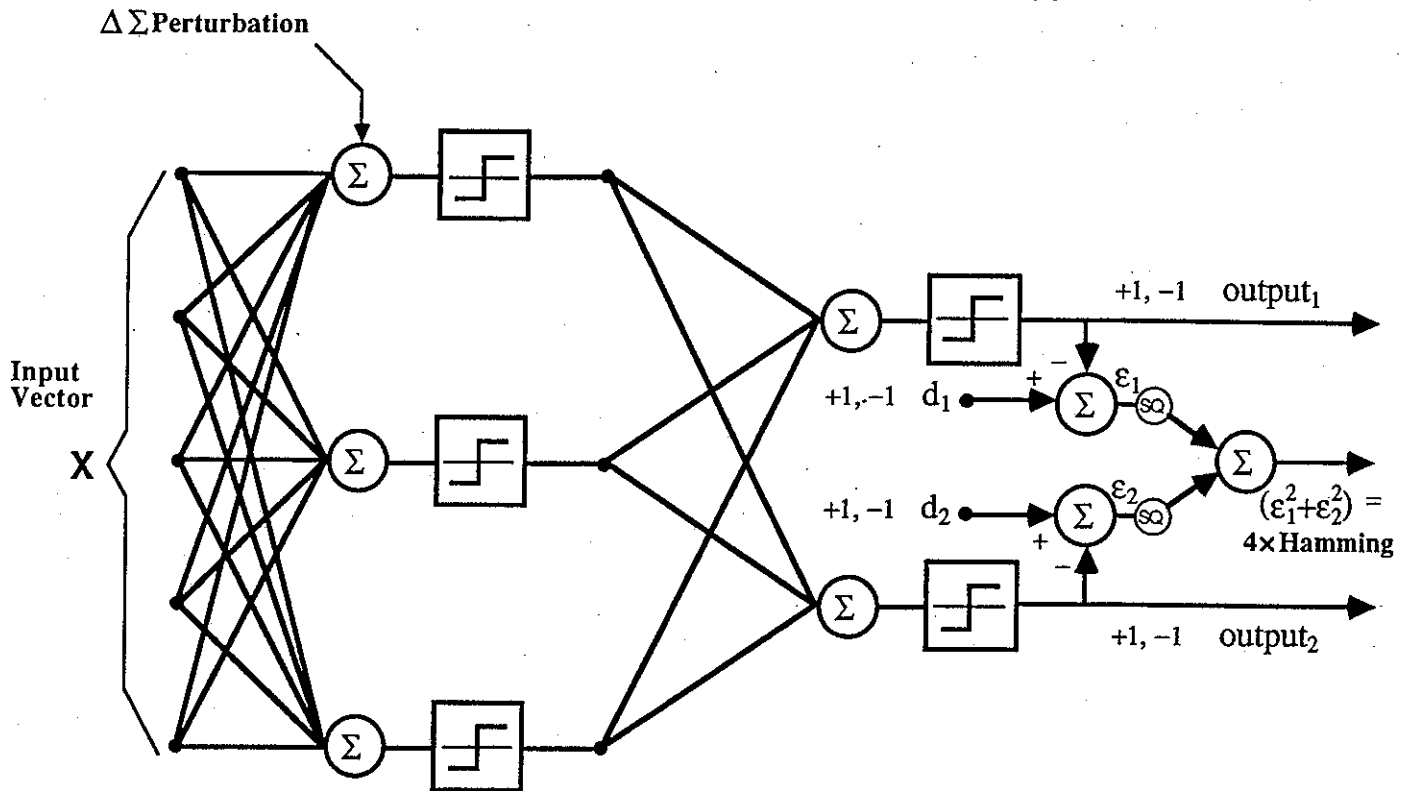
- Performs poorly if training set is not linearly separable.





**A five-Adaline example of the MRI Architecture**

## MRLI of B. Widrow & R. Winter



For each layer, beginning with layer 1:

Toggle output of neuron with sum closest to zero. If output Hamming error is reduced, adapt neuron. Repeat for neuron whose sum is next closest to zero, etc. Can also adapt two at a time, etc. Adaptation reduces Hamming error.