# Example

# Elementary Decision Boundaries
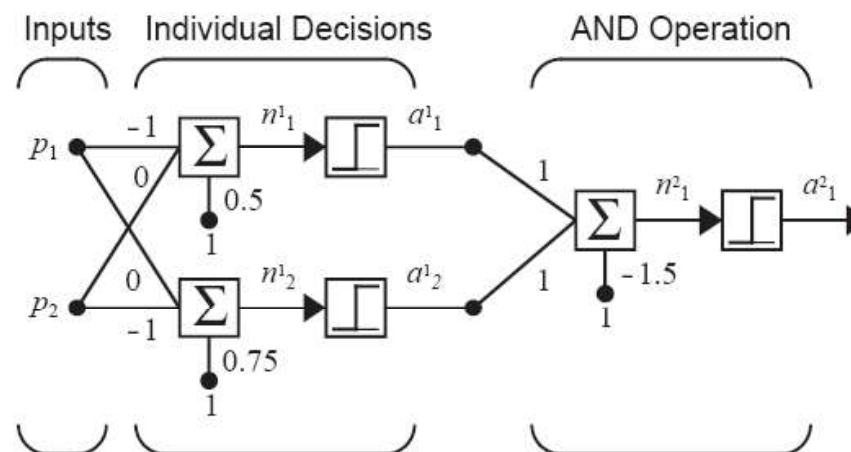


First Boundary:

$$a_1^1 = hardlim(\begin{bmatrix} -1 & 0 \end{bmatrix}\mathbf{p} + 0.5)$$

Second Boundary:

$$a_2^1 = hardlim(\begin{bmatrix} 0 & -1 \end{bmatrix}\mathbf{p} + 0.75)$$

First Subnetwork

Inputs  Individual Decisions  AND Operation

# Elementary Decision Boundaries



Third Boundary:

$$a_3^1 = hardlim([1 \; 0]\mathbf{p} - 1.5)$$

Fourth Boundary:

$$a_4^1 = hardlim([0 \; 1]\mathbf{p} - 0.25)$$
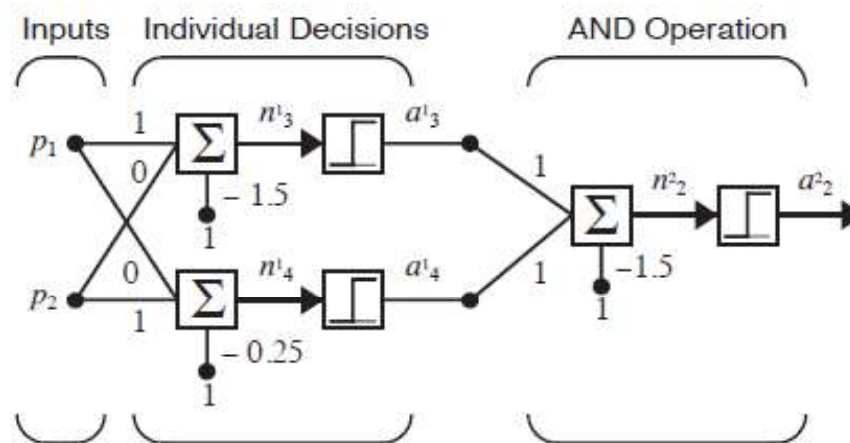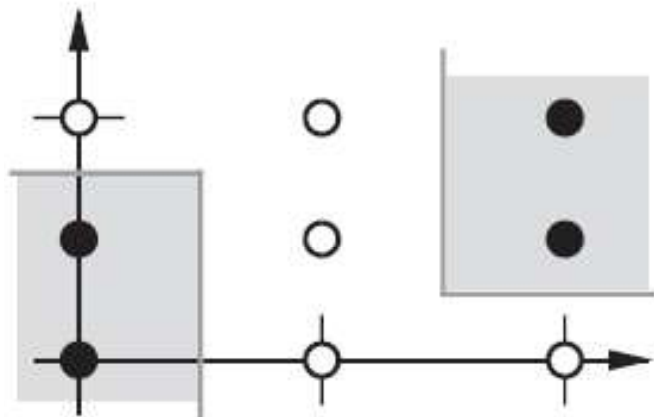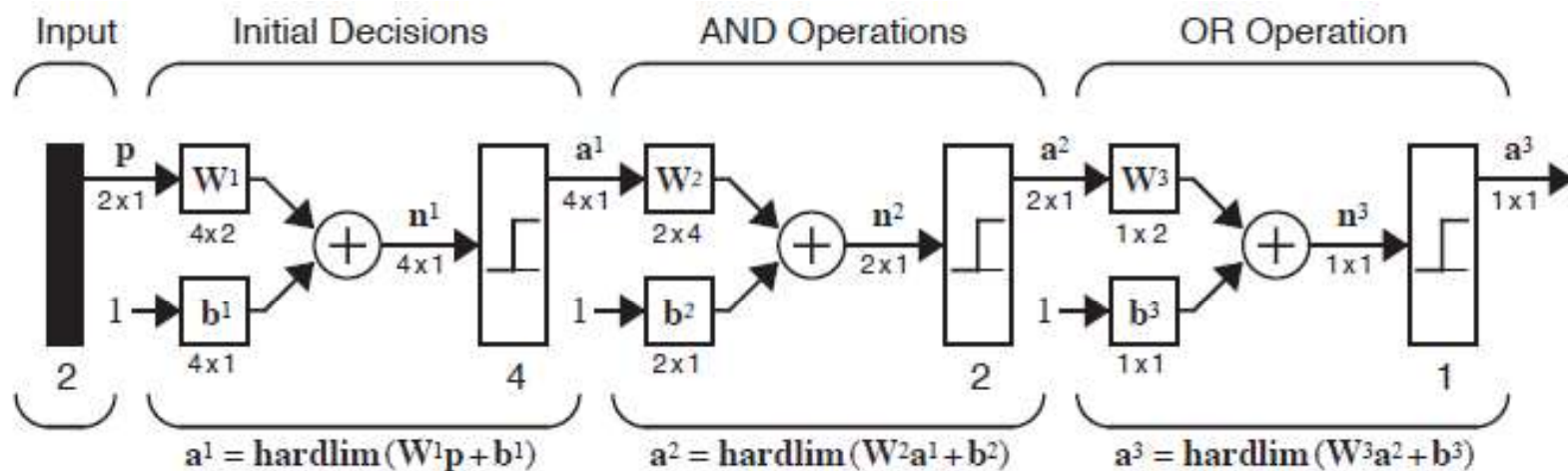
Second Subnetwork

# Total Network

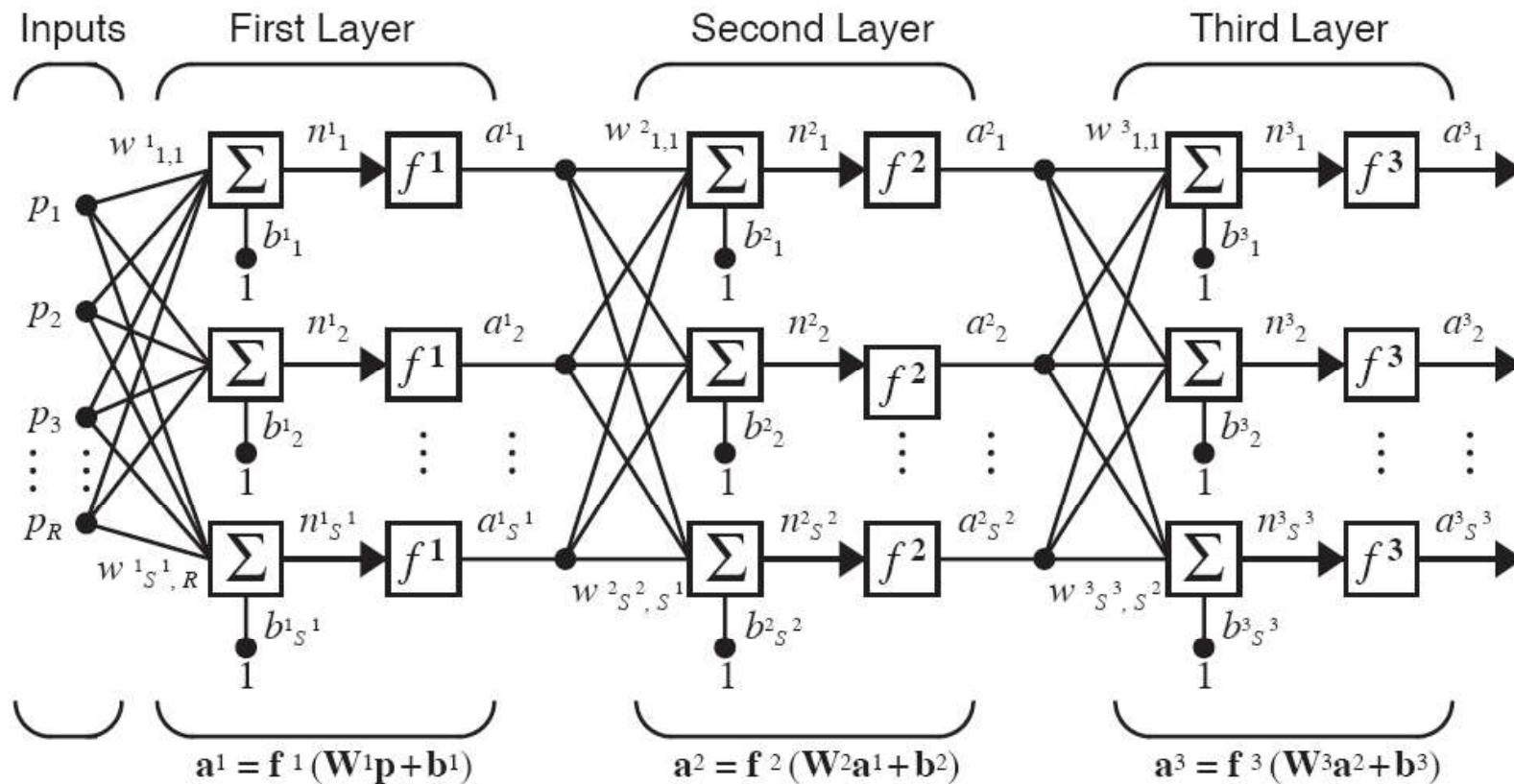$$\mathbf{W}^1 = \begin{bmatrix} -1 & 0 \\ 0 & -1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \qquad \mathbf{b}^1 = \begin{bmatrix} 0.5 \\ 0.75 \\ -1.5 \\ -0.25 \end{bmatrix}$$

$$\mathbf{W}^2 = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \qquad \mathbf{b}^2 = \begin{bmatrix} -1.5 \\ -1.5 \end{bmatrix}$$

$$\mathbf{W}^3 = \begin{bmatrix} 1 & 1 \end{bmatrix} \qquad \mathbf{b}^3 = \begin{bmatrix} -0.5 \end{bmatrix}$$

| Input | Initial Decisions | AND Operations | OR Operation |
|---|---|---|---|

**p** 2x1 → **W¹** 4x2 ; 1 → **b¹** 4x1 ; ⊕ → **n¹** 4x1 → (hardlim) → **a¹** 4x1
**W²** 2x4 ; 1 → **b²** 2x1 ; ⊕ → **n²** 2x1 → (hardlim) → **a²** 2x1
**W³** 1x2 ; 1 → **b³** 1x1 ; ⊕ → **n³** 1x1 → (hardlim) → **a³** 1x1

2

4

2

1

$$a^1 = \mathrm{hardlim}(\mathbf{W}^1 p + b^1) \qquad a^2 = \mathrm{hardlim}(\mathbf{W}^2 a^1 + b^2) \qquad a^3 = \mathrm{hardlim}(\mathbf{W}^3 a^2 + b^3)$$

5

# Multilayer Perceptron



Inputs — First Layer — Second Layer — Third Layer

$$\mathbf{a}^1 = \mathbf{f}^1(\mathbf{W}^1\mathbf{p}+\mathbf{b}^1)$$

$$\mathbf{a}^2 = \mathbf{f}^2(\mathbf{W}^2\mathbf{a}^1+\mathbf{b}^2)$$

$$\mathbf{a}^3 = \mathbf{f}^3(\mathbf{W}^3\mathbf{a}^2+\mathbf{b}^3)$$
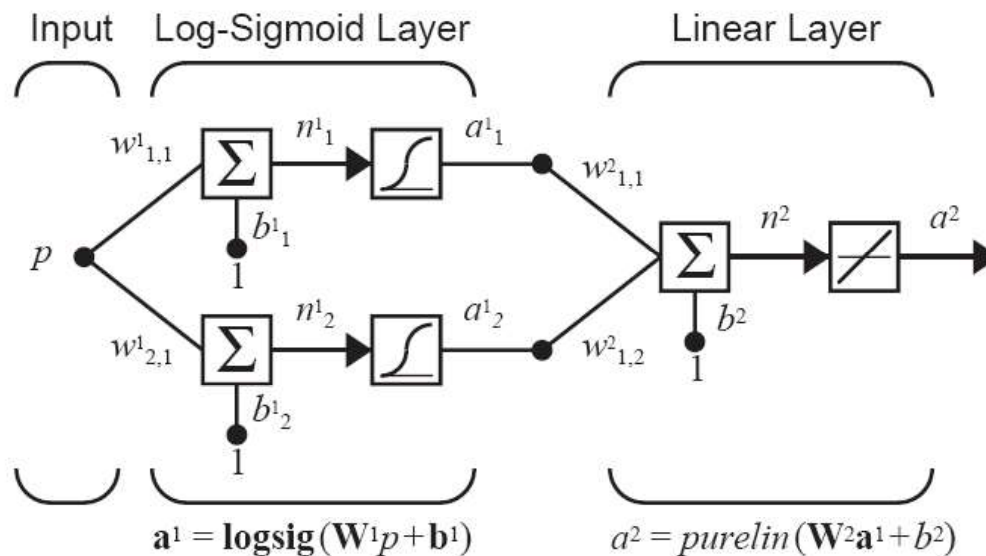
$$\mathbf{a}^3 = \mathbf{f}^3(\mathbf{W}^3\mathbf{f}^2(\mathbf{W}^2\mathbf{f}^1(\mathbf{W}^1\mathbf{p}+\mathbf{b}^1)+\mathbf{b}^2)+\mathbf{b}^3)$$

$$R - S^1 - S^2 - S^3 \text{ Network}$$

# Function Approximation Example



Input  Log-Sigmoid Layer          Linear Layer

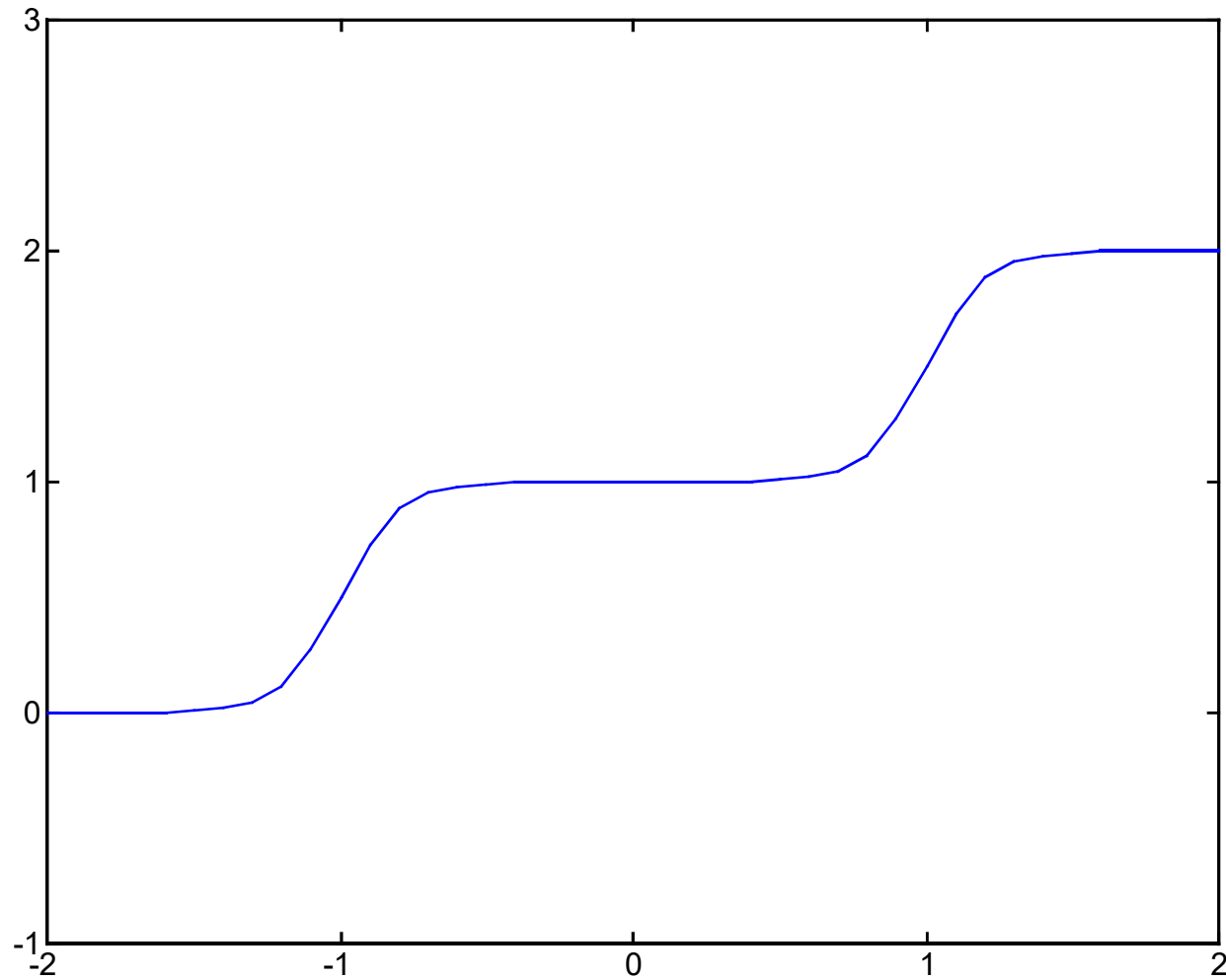$$a^1 = logsig(W^1 p + b^1)$$ $$a^2 = purelin(W^2 a^1 + b^2)$$

$$f^1(n) = \frac{1}{1 + e^{-n}}$$

$$f^2(n) = n$$
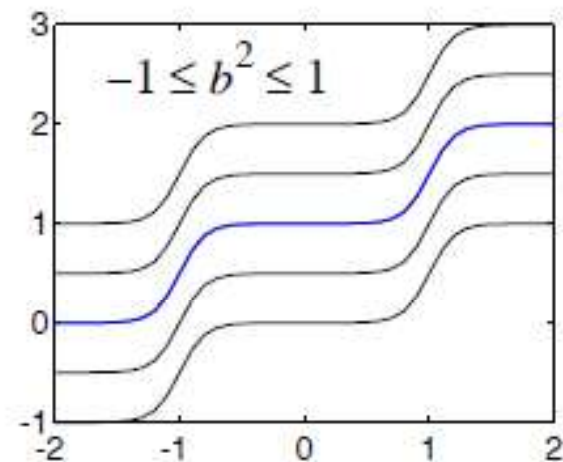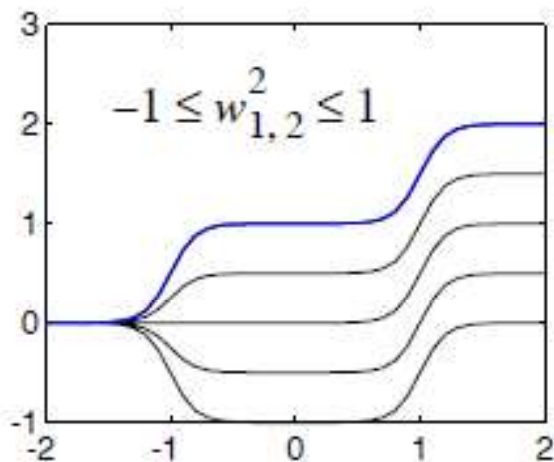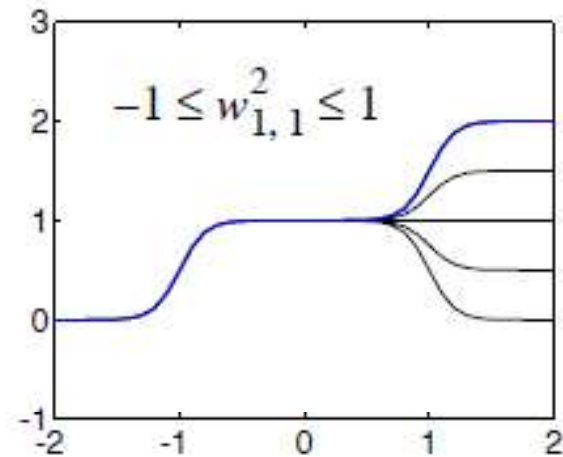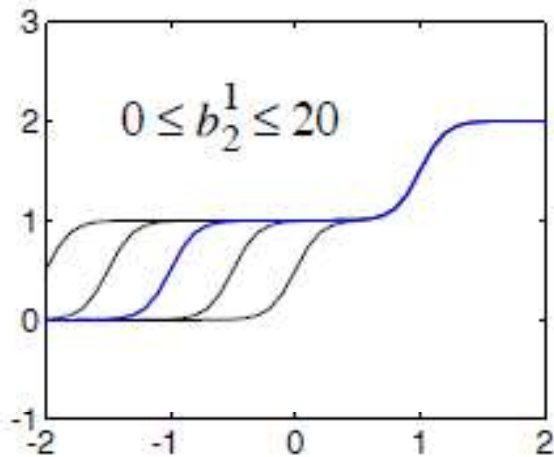
## Nominal Parameter Values

$$w^1_{1,1} = 10 \qquad w^1_{2,1} = 10 \qquad b^1_1 = -10 \qquad b^1_2 = 10$$

$$w^2_{1,1} = 1 \qquad w^2_{1,2} = 1 \qquad b^2 = 0$$

# Nominal Response

# Parameter Variations

# Multilayer Network



Hidden Layers

Output Layer

Input    First Layer    Second Layer    Third Layer

$$\mathbf{a}^1 = \mathbf{f}^1(\mathbf{W}^1\mathbf{p} + \mathbf{b}^1) \qquad \mathbf{a}^2 = \mathbf{f}^2(\mathbf{W}^2\mathbf{a}^1 + \mathbf{b}^2) \qquad \mathbf{a}^3 = \mathbf{f}^3(\mathbf{W}^3\mathbf{a}^2 + \mathbf{b}^3)$$

$$\mathbf{a}^3 = \mathbf{f}^3(\mathbf{W}^3\mathbf{f}^2(\mathbf{W}^2\mathbf{f}^1(\mathbf{W}^1\mathbf{p} + \mathbf{b}^1) + \mathbf{b}^2) + \mathbf{b}^3)$$

10
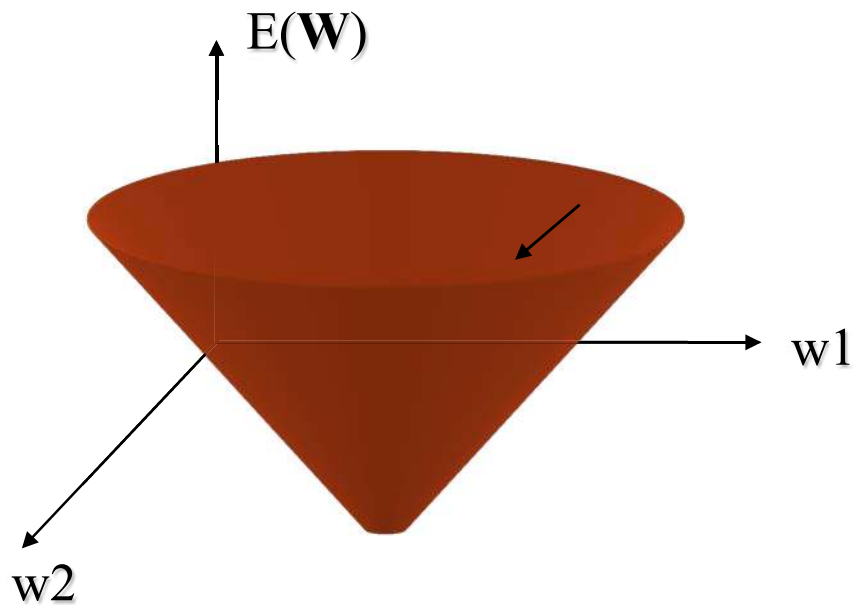
# الگوریتم Gradient descent

با توجه به نحوه تعریف E ، سطح خطا به صورت یک سهمی خواهد بود. ما به دنبال وزن‌هایی هستیم که حداقل خطا را داشته باشند . الگوریتم Gradient descent در فضای وزن‌ها به دنبال برداری می‌گردد که خطا را حداقل کند.

E(**W**)

w1

w2

این الگوریتم از یک مقدار دلخواه برای بردار وزن شروع کرده و در هر مرحله وزن‌ها را طوری تغییر می‌دهد که در جهت شیب کاهشی منحنی فوق خطا کاهش داده شود.

# Perceptron Vs. LMS

**Derivative**

- **Linear activation function has derivative**

  but

- **Sign (bipolar, unipolar) has not derivative**

# Chain Rule

$$\frac{df(n(w))}{dw} = \frac{df(n)}{dn} \times \frac{dn(w)}{dw}$$

Example

$$f(n) = \cos(n) \qquad n = e^{2w} \qquad f(n(w)) = \cos(e^{2w})$$

$$\frac{df(n(w))}{dw} = \frac{df(n)}{dn} \times \frac{dn(w)}{dw} = (-\sin(n))(2e^{2w}) = (-\sin(e^{2w}))(2e^{2w})$$

# Performance Index

Training Set

$$\{\mathbf{p}_1, \mathbf{t}_1\}, \{\mathbf{p}_2, \mathbf{t}_2\}, \dots, \{\mathbf{p}_Q, \mathbf{t}_Q\}$$

Mean Square Errors

$$F(\mathbf{x}) = E[e^2] = E[(t-a)^2]$$

Vector Case

$$F(\mathbf{x}) = E[\mathbf{e}^T \mathbf{e}] = E[(\mathbf{t}-\mathbf{a})^T(\mathbf{t}-\mathbf{a})]$$
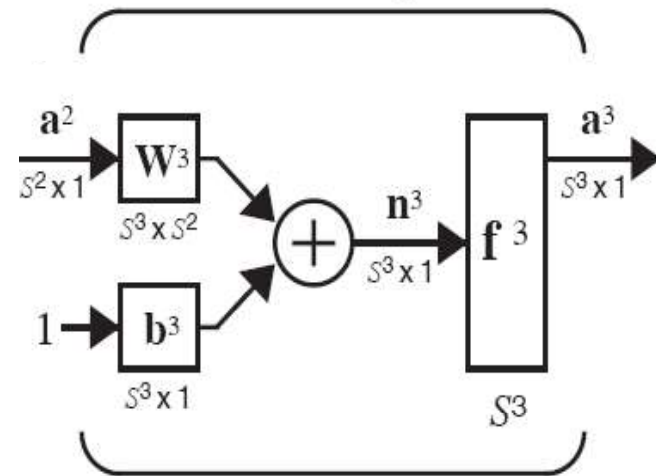
Approximate Mean Square Error (Single Sample)

$$\hat{F}(\mathbf{x}) = (\mathbf{t}(k) - \mathbf{a}(k))^T(\mathbf{t}(k) - \mathbf{a}(k)) = \mathbf{e}^T(k)\mathbf{e}(k)$$

14

# A Simple Example

$\mathbf{n}^3 = (\mathbf{W}^3\mathbf{a}^2 + \mathbf{b}^3)$

$\mathbf{a}^3 = \mathbf{f}^3(\mathbf{W}^3\mathbf{a}^2 + \mathbf{b}^3)$

$\mathbf{w}^m(k+1) = \mathbf{w}^m(k) - \alpha\dfrac{\partial\hat{F}}{\partial\mathbf{w}^m}$



$$\frac{\partial\hat{F}}{\partial W^3} = \frac{\partial\hat{F}}{\partial e}\cdot\frac{\partial e}{\partial a^3}\cdot\frac{\partial a^3}{\partial n^3}\cdot\frac{\partial n^3}{\partial W^3}$$

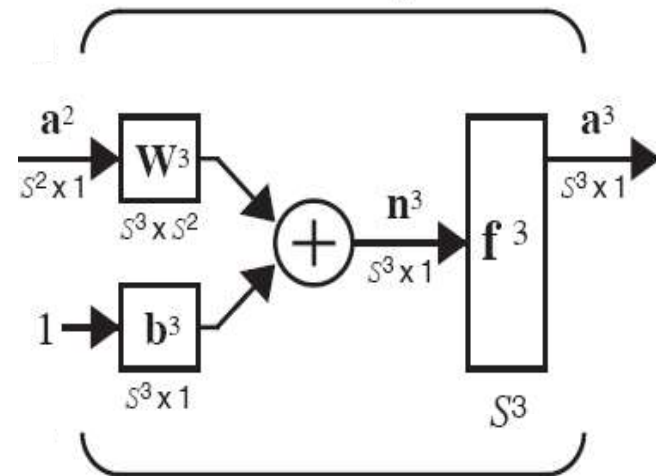$\dfrac{\partial\hat{F}}{\partial e} = 2e$ $\qquad$ $\dfrac{\partial e}{\partial a^3} = -1$ $\qquad$ $\dfrac{\partial a^3}{\partial n^3} = \dot{f}^3(n^3)$ $\qquad$ $\dfrac{\partial n^3}{\partial W^3} = a^2$

# A Simple Example

$$\mathbf{n}^3 = (\mathbf{W}^3\mathbf{a}^2 + \mathbf{b}^3)$$

$$\mathbf{a}^3 = \mathbf{f}^3(\mathbf{W}^3\mathbf{a}^2 + \mathbf{b}^3)$$

$$b^m(k+1) = b^m(k) - \alpha\frac{\partial\hat{F}}{\partial b^m}$$



$$\frac{\partial\hat{F}}{\partial b^3} = \frac{\partial\hat{F}}{\partial e}\cdot\frac{\partial e}{\partial a^3}\cdot\frac{\partial a^3}{\partial n^3}\cdot\frac{\partial n^3}{\partial b^3}$$

$$\frac{\partial\hat{F}}{\partial e} = 2e \qquad \frac{\partial e}{\partial a^3} = -1 \qquad \frac{\partial a^3}{\partial n^3} = \dot{f}^3(n^3) \qquad \frac{\partial n^3}{\partial b^3} = 1$$

16

# Chain Rule

## Approximate Steepest Descent

$$w_{i,j}^m(k+1) = w_{i,j}^m(k) - \alpha \frac{\partial \hat{F}}{\partial w_{i,j}^m} \qquad b_i^m(k+1) = b_i^m(k) - \alpha \frac{\partial \hat{F}}{\partial b_i^m}$$

## Application to Gradient Calculation

$$\frac{\partial \hat{F}}{\partial w_{i,j}^m} = \frac{\partial \hat{F}}{\partial n_i^m} \times \frac{\partial n_i^m}{\partial w_{i,j}^m} \qquad\qquad \frac{\partial \hat{F}}{\partial b_i^m} = \frac{\partial \hat{F}}{\partial n_i^m} \times \frac{\partial n_i^m}{\partial b_i^m}$$

17

# Gradient Calculation

$$n_i^m = \sum_{j=1}^{S^{m-1}} w_{i,j}^m a_j^{m-1} + b_i^m$$

$$\frac{\partial n_i^m}{\partial w_{i,j}^m} = a_j^{m-1} \qquad \frac{\partial n_i^m}{\partial b_i^m} = 1$$

## Sensitivity

$$s_i^m \equiv \frac{\partial \hat{F}}{\partial n_i^m}$$

## Gradient

$$\frac{\partial \hat{F}}{\partial w_{i,j}^m} = s_i^m a_j^{m-1} \qquad \frac{\partial \hat{F}}{\partial b_i^m} = s_i^m$$

# Steepest Descent

$$w_{i,j}^m(k+1) = w_{i,j}^m(k) - \alpha s_i^m a_j^{m-1} \qquad b_i^m(k+1) = b_i^m(k) - \alpha s_i^m$$

$$\mathbf{W}^m(k+1) = \mathbf{W}^m(k) - \alpha \mathbf{s}^m (\mathbf{a}^{m-1})^T \qquad \mathbf{b}^m(k+1) = \mathbf{b}^m(k) - \alpha \mathbf{s}^m$$

$$\mathbf{s}^m \equiv \frac{\partial \hat{F}}{\partial \mathbf{n}^m} = \begin{bmatrix} \dfrac{\partial \hat{F}}{\partial n_1^m} \\[2ex] \dfrac{\partial \hat{F}}{\partial n_2^m} \\[1ex] \vdots \\[1ex] \dfrac{\partial \hat{F}}{\partial n_{S^m}^m} \end{bmatrix}$$

Next Step: Compute the Sensitivities (Backpropagation)

19

# Jacobian Matrix

$$\frac{\partial \mathbf{n}^{m+1}}{\partial \mathbf{n}^m} \equiv \begin{bmatrix} \dfrac{\partial n_1^{m+1}}{\partial n_1^m} & \dfrac{\partial n_1^{m+1}}{\partial n_2^m} & \cdots & \dfrac{\partial n_1^{m+1}}{\partial n_{S^m}^m} \\[2ex] \dfrac{\partial n_2^{m+1}}{\partial n_1^m} & \dfrac{\partial n_2^{m+1}}{\partial n_2^m} & \cdots & \dfrac{\partial n_2^{m+1}}{\partial n_{S^m}^m} \\[2ex] \vdots & \vdots & & \vdots \\[2ex] \dfrac{\partial n_{S^{m+1}}^{m+1}}{\partial n_1^m} & \dfrac{\partial n_{S^{m+1}}^{m+1}}{\partial n_2^m} & \cdots & \dfrac{\partial n_{S^{m+1}}^{m+1}}{\partial n_{S^m}^m} \end{bmatrix}$$

$$\frac{\partial n_i^{m+1}}{\partial n_j^m} = \frac{\partial \left( \sum\limits_{l=1}^{S^m} w_{i,l}^{m+1} a_l^m + b_i^{m+1} \right)}{\partial n_j^m} = w_{i,j}^{m+1} \frac{\partial a_j^m}{\partial n_j^m}$$

$$\frac{\partial n_i^{m+1}}{\partial n_j^m} = w_{i,j}^{m+1} \frac{\partial f^m(n_j^m)}{\partial n_j^m} = w_{i,j}^{m+1} \dot{f}^m(n_j^m)$$

$$\dot{f}^m(n_j^m) = \frac{\partial f^m(n_j^m)}{\partial n_j^m}$$

$$\frac{\partial \mathbf{n}^{m+1}}{\partial \mathbf{n}^m} = \mathbf{W}^{m+1} \dot{\mathbf{F}}^m(\mathbf{n}^m)$$

$$\dot{\mathbf{F}}^m(\mathbf{n}^m) = \begin{bmatrix} \dot{f}^m(n_1^m) & 0 & \cdots & 0 \\ 0 & \dot{f}^m(n_2^m) & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & \dot{f}^m(n_{S^m}^m) \end{bmatrix}$$

# Backpropagation (Sensitivities)

$$s^m = \frac{\partial \hat{F}}{\partial \mathbf{n}^m} = \left(\frac{\partial \mathbf{n}^{m+1}}{\partial \mathbf{n}^m}\right)^T \frac{\partial \hat{F}}{\partial \mathbf{n}^{m+1}} = \dot{\mathbf{F}}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T \frac{\partial \hat{F}}{\partial \mathbf{n}^{m+1}}$$

$$\boxed{s^m = \dot{\mathbf{F}}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T s^{m+1}}$$

The sensitivities are computed by starting at the last layer, and then propagating backwards through the network to the first layer.

$$s^M \rightarrow s^{M-1} \rightarrow \ldots \rightarrow s^2 \rightarrow s^1$$

21

# Initialization (Last Layer)

$$s_i^M = \frac{\partial \hat{F}}{\partial n_i^M} = \frac{\partial(\mathbf{t}-\mathbf{a})^T(\mathbf{t}-\mathbf{a})}{\partial n_i^M} = \frac{\partial \sum\limits_{j=1}^{S^M} (t_j - a_j)^2}{\partial n_i^M} = -2(t_i - a_i)\frac{\partial a_i}{\partial n_i^M}$$

$$\frac{\partial a_i}{\partial n_i^M} = \frac{\partial a_i^M}{\partial n_i^M} = \frac{\partial f^M(n_i^M)}{\partial n_i^M} = \dot{f}^M(n_i^M)$$

$$s_i^M = -2(t_i - a_i)\dot{f}^M(n_i^M)$$

$$s^M = -2\dot{\mathbf{F}}^M(\mathbf{n}^M)(\mathbf{t}-\mathbf{a})$$

# Summary

## Forward Propagation

$$\mathbf{a}^0 = \mathbf{p}$$

$$\mathbf{a}^{m+1} = \mathbf{f}^{m+1}(\mathbf{W}^{m+1}\mathbf{a}^m + \mathbf{b}^{m+1}) \qquad m = 0, 2, \ldots, M-1$$

$$\mathbf{a} = \mathbf{a}^M$$

## Backpropagation

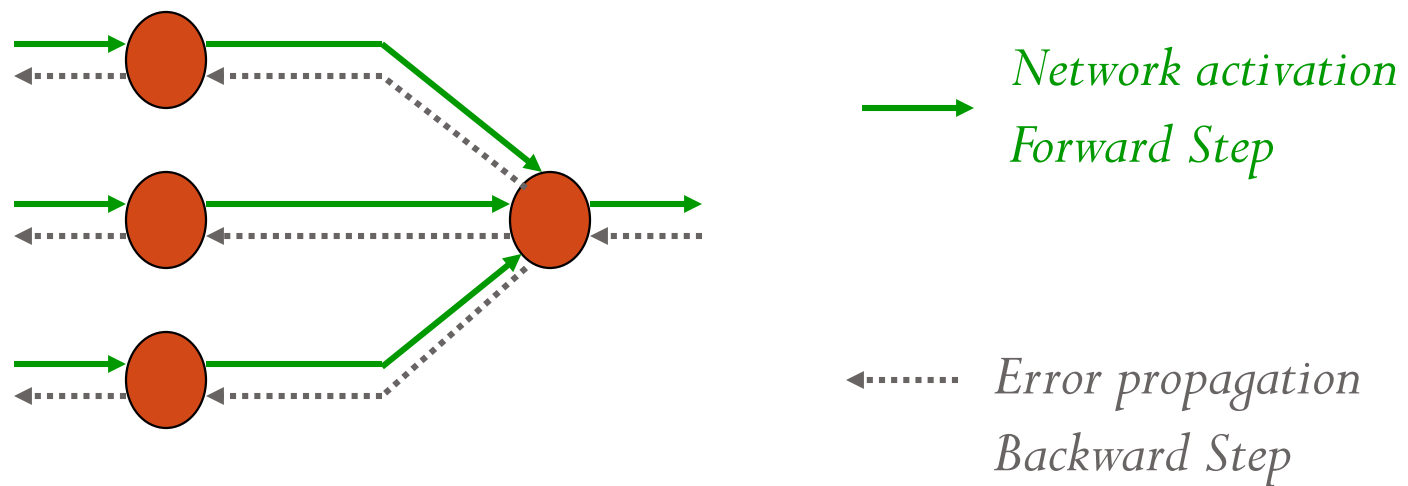$$\mathbf{s}^M = -2\dot{\mathbf{F}}^M(\mathbf{n}^M)(\mathbf{t} - \mathbf{a})$$

$$\mathbf{s}^m = \dot{\mathbf{F}}^m(\mathbf{n}^m)(\mathbf{W}^{m+1})^T\mathbf{s}^{m+1} \qquad m = M-1, \ldots, 2, 1$$

## Weight Update

$$\mathbf{W}^m(k+1) = \mathbf{W}^m(k) - \alpha\mathbf{s}^m(\mathbf{a}^{m-1})^T \qquad \mathbf{b}^m(k+1) = \mathbf{b}^m(k) - \alpha\mathbf{s}^m$$
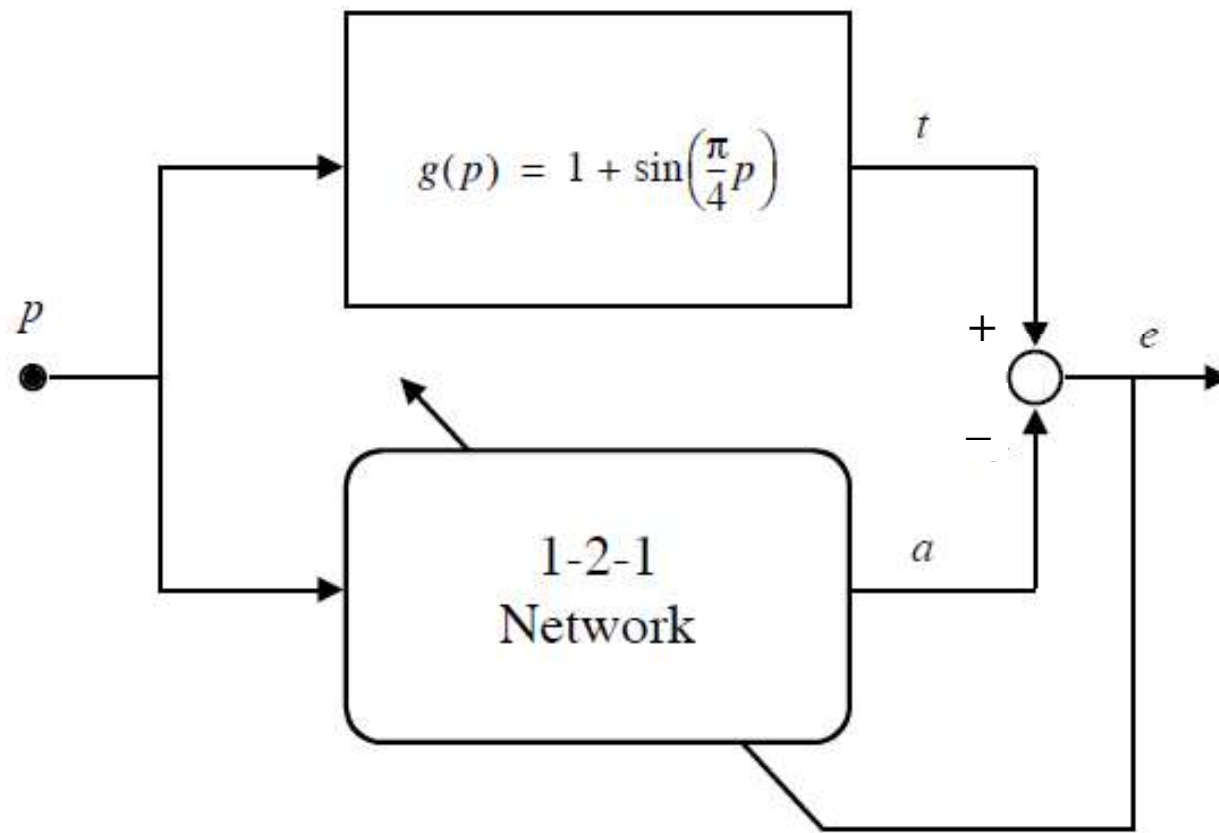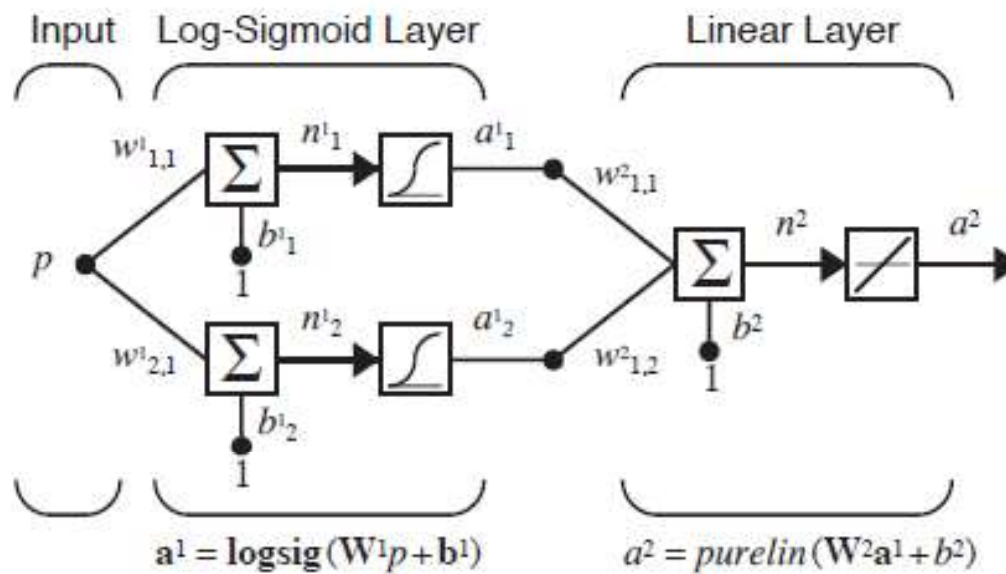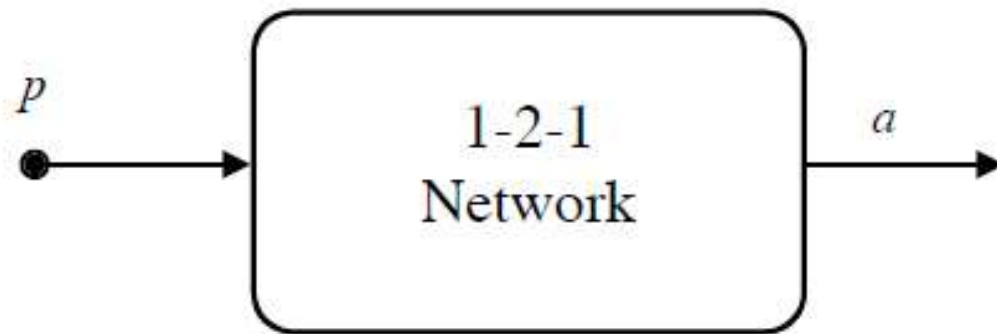
23

# Summary

- Back-propagation training algorithm



*Network activation*
*Forward Step*

*Error propagation*
*Backward Step*

- Backpropagation adjusts the weights of the NN in order to minimize the network total mean squared error.
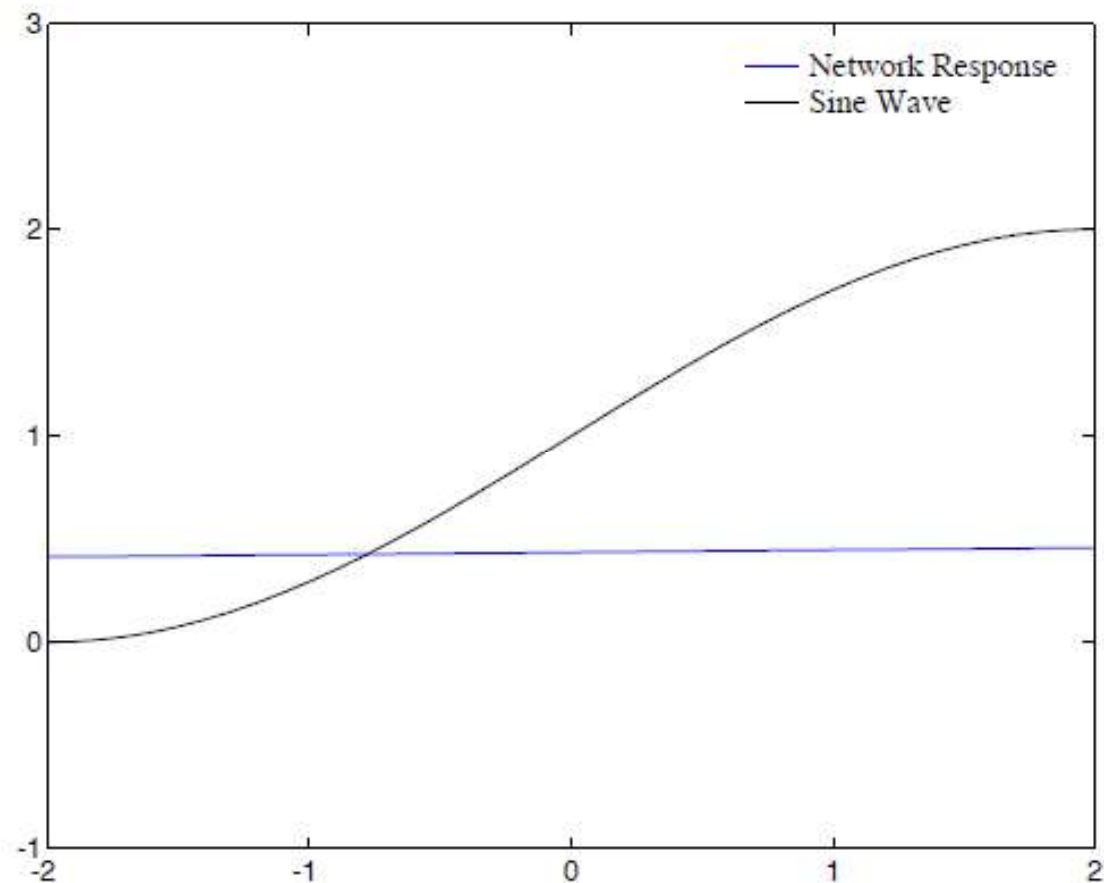
# Example: Function Approximation



$$g(p) = 1 + \sin\left(\frac{\pi}{4}p\right)$$

1-2-1 Network

# Network



Input    Log-Sigmoid Layer        Linear Layer

$$\mathbf{a}^1 = logsig\,(\mathbf{W}^1 p + \mathbf{b}^1) \qquad a^2 = purelin\,(\mathbf{W}^2 \mathbf{a}^1 + b^2)$$

# Initial Conditions

$$\mathbf{W}^1(0) = \begin{bmatrix} -0.27 \\ -0.41 \end{bmatrix} \quad \mathbf{b}^1(0) = \begin{bmatrix} -0.48 \\ -0.13 \end{bmatrix} \quad \mathbf{W}^2(0) = \begin{bmatrix} 0.09 & -0.17 \end{bmatrix} \quad \mathbf{b}^2(0) = \begin{bmatrix} 0.48 \end{bmatrix}$$

# Forward Propagation

$$a^0 = p = 1$$

$$\mathbf{a}^1 = \mathbf{f}^1(\mathbf{W}^1\mathbf{a}^0 + \mathbf{b}^1) = \text{logsig}\left(\begin{bmatrix} -0.27 \\ -0.41 \end{bmatrix}\begin{bmatrix} 1 \end{bmatrix} + \begin{bmatrix} -0.48 \\ -0.13 \end{bmatrix}\right) = \text{logsig}\left(\begin{bmatrix} -0.75 \\ -0.54 \end{bmatrix}\right)$$

$$\mathbf{a}^1 = \begin{bmatrix} \dfrac{1}{1 + e^{0.75}} \\[2ex] \dfrac{1}{1 + e^{0.54}} \end{bmatrix} = \begin{bmatrix} 0.321 \\ 0.368 \end{bmatrix}$$

$$a^2 = \mathbf{f}^2(\mathbf{W}^2\mathbf{a}^1 + \mathbf{b}^2) = purelin\left(\begin{bmatrix} 0.09 & -0.17 \end{bmatrix}\begin{bmatrix} 0.321 \\ 0.368 \end{bmatrix} + \begin{bmatrix} 0.48 \end{bmatrix}\right) = \begin{bmatrix} 0.446 \end{bmatrix}$$

$$e = t - a = \left\{1 + \sin\left(\frac{\pi}{4}p\right)\right\} - a^2 = \left\{1 + \sin\left(\frac{\pi}{4}1\right)\right\} - 0.446 = 1.261$$

# Transfer Function Derivatives

$$\dot{f}^1(n) = \frac{d}{dn}\left(\frac{1}{1+e^{-n}}\right) = \frac{e^{-n}}{(1+e^{-n})^2} = \left(1 - \frac{1}{1+e^{-n}}\right)\left(\frac{1}{1+e^{-n}}\right) = (1-a^1)a^1$$

$$\dot{f}^2(n) = \frac{d}{dn}(n) = 1$$

# Backpropagation

$$\mathbf{s}^2 = -2\dot{\mathbf{F}}^2(\mathbf{n}^2)(\mathbf{t} - \mathbf{a}) = -2\left[\dot{f}^2(n^2)\right](1.261) = -2\left[1\right](1.261) = -2.522$$

$$\mathbf{s}^1 = \dot{\mathbf{F}}^1(\mathbf{n}^1)(\mathbf{W}^2)^T\mathbf{s}^2 = \begin{bmatrix} (1-a_1^1)(a_1^1) & 0 \\ 0 & (1-a_2^1)(a_2^1) \end{bmatrix} \begin{bmatrix} 0.09 \\ -0.17 \end{bmatrix} \begin{bmatrix} -2.522 \end{bmatrix}$$

$$\mathbf{s}^1 = \begin{bmatrix} (1-0.321)(0.321) & 0 \\ 0 & (1-0.368)(0.368) \end{bmatrix} \begin{bmatrix} 0.09 \\ -0.17 \end{bmatrix} \begin{bmatrix} -2.522 \end{bmatrix}$$

$$\mathbf{s}^1 = \begin{bmatrix} 0.218 & 0 \\ 0 & 0.233 \end{bmatrix} \begin{bmatrix} -0.227 \\ 0.429 \end{bmatrix} = \begin{bmatrix} -0.0495 \\ 0.0997 \end{bmatrix}$$
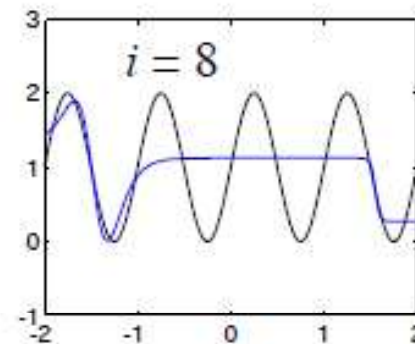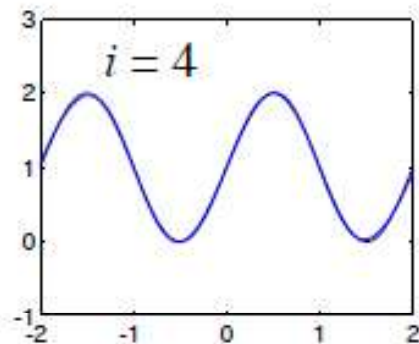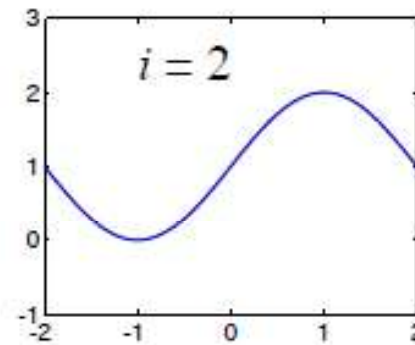
# Weight Update

$$\alpha = 0.1$$

$$\mathbf{W}^2(1) = \mathbf{W}^2(0) - \alpha \mathbf{s}^2(\mathbf{a}^1)^T = \begin{bmatrix} 0.09 & -0.17 \end{bmatrix} - 0.1 \begin{bmatrix} -2.522 \end{bmatrix} \begin{bmatrix} 0.321 & 0.368 \end{bmatrix}$$

$$\mathbf{W}^2(1) = \begin{bmatrix} 0.171 & -0.0772 \end{bmatrix}$$

$$\mathbf{b}^2(1) = \mathbf{b}^2(0) - \alpha \mathbf{s}^2 = \begin{bmatrix} 0.48 \end{bmatrix} - 0.1 \begin{bmatrix} -2.522 \end{bmatrix} = \begin{bmatrix} 0.732 \end{bmatrix}$$
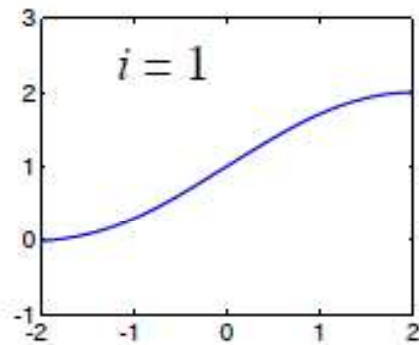
$$\mathbf{W}^1(1) = \mathbf{W}^1(0) - \alpha \mathbf{s}^1(\mathbf{a}^0)^T = \begin{bmatrix} -0.27 \\ -0.41 \end{bmatrix} - 0.1 \begin{bmatrix} -0.0495 \\ 0.0997 \end{bmatrix} \begin{bmatrix} 1 \end{bmatrix} = \begin{bmatrix} -0.265 \\ -0.420 \end{bmatrix}$$

$$\mathbf{b}^1(1) = \mathbf{b}^1(0) - \alpha \mathbf{s}^1 = \begin{bmatrix} -0.48 \\ -0.13 \end{bmatrix} - 0.1 \begin{bmatrix} -0.0495 \\ 0.0997 \end{bmatrix} = \begin{bmatrix} -0.475 \\ -0.140 \end{bmatrix}$$
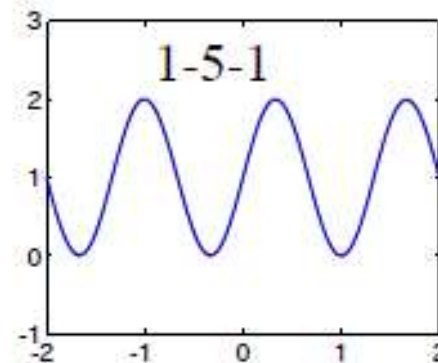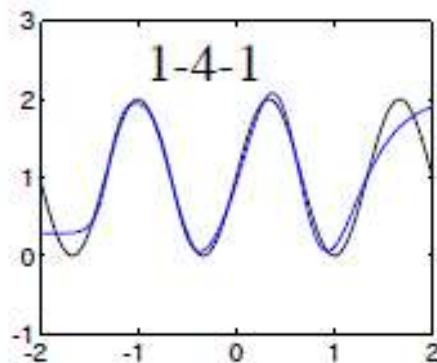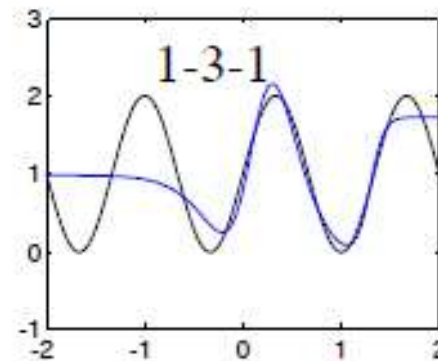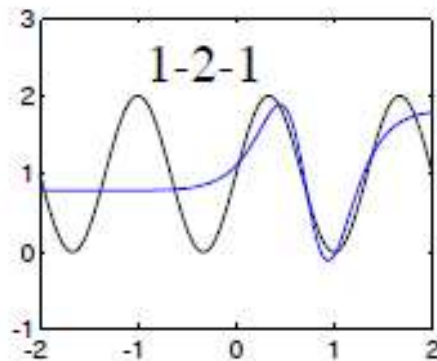
31

# Choice of Network Architecture

$$g(p) = 1 + \sin\left(\frac{i\pi}{4}p\right) \text{ for } -2 \le p \le 2$$

1-3-1 Network

# Choice of Network Architecture

$$g(p) = 1 + \sin\left(\frac{6\pi}{4}p\right) \text{ for } -2 \leq p \leq 2$$
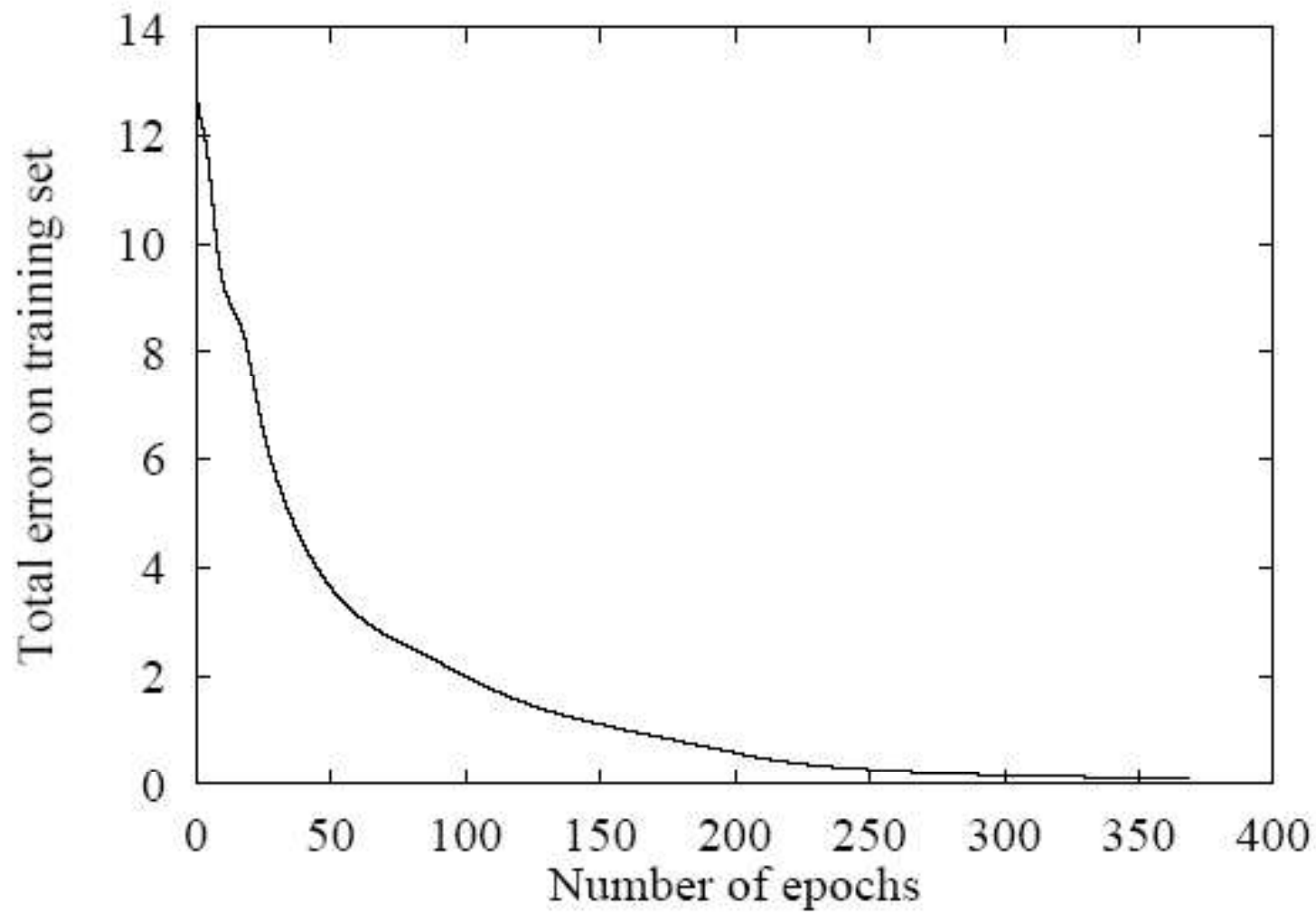
# Disadvantage of BP algorithm

- Slow convergence speed

- Sensitivity to initial conditions

- Trapped in local minima

- Instability if learning rate is too large


- Note: despite above disadvantages, it is popularly used in control community. There are numerous extensions to improve BP algorithm.

# شرط خاتمه

معمولا الگوریتم BP پیش از خاتمه هزاران بار با استفاده از همان دادههای آموزشی تکرار میشود. در این حالت شروط مختلفی را میتوان برای خاتمه الگوریتم به کار برد:

- توقف بعد از تکرار به دفعات معین،

- توقف وقتی که خطا از یک مقدار تعیین شده کمتر شود،

- توقف وقتی که خطا در مثالهای مجموعه اعتبار سنجی از قاعده خاصی پیروی کند.

اگر دفعات تکرار کم باشد خطا خواهیم داشت و اگر زیاد باشد مساله Overfitting رخ خواهد داد.

# منحنی یادگیری

# قدرت تعمیم و Overfitting

- شرط پایان الگوریتم BP چیست؟

- یک انتخاب این است که الگوریتم را آنقدر ادامه دهیم تا خطا از مقدار معینی کمتر شود. این امر می‌تواند منجر به Overfitting شود.



Number of Epochs