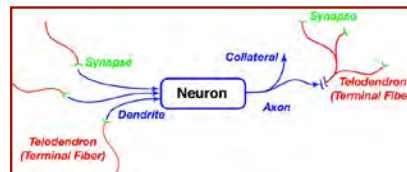


Introduction to Neural Networks

Robert Stengel

Robotics and Intelligent Systems, MAE 345,
Princeton University, 2017

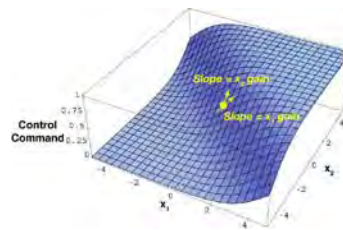
- Natural and artificial neurons
- Natural and computational neural networks
 - Linear network
 - Perceptron
 - Sigmoid network
 - Radial basis function
- Applications of neural networks
- Supervised training
 - Left pseudoinverse
 - Steepest descent
 - Back-propagation



Copyright 2017 by Robert Stengel. All rights reserved. For educational use only.
<http://www.princeton.edu/~stengel/MAE345.html>

1

Applications of Computational Neural Networks

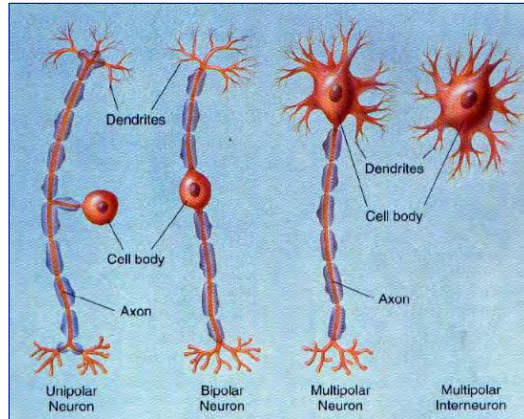


- Classification of data sets
- Image processing
- Language interpretation
- Nonlinear function approximation
 - Efficient data storage and retrieval
 - System identification
- Nonlinear and adaptive control systems

2

Neurons

- **Biological cells with significant electrochemical activity**
- **~10-100 billion neurons in the brain**
- **Inputs from thousands of other neurons**
- **Output is scalar, but may have thousands of branches**

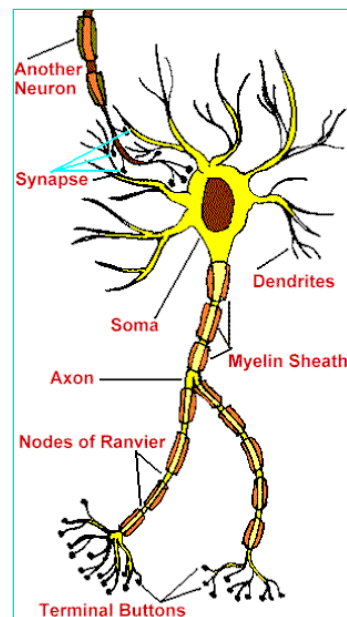


- **Afferent (sensor) neurons** send signals from organs and the periphery to the central nervous system
- **Efferent (motor) neurons** issue commands from the CNS to effector (e.g., muscle) cells
- **Interneurons** send signals between neurons in the central nervous system
- Signals are **ionic**, i.e., chemical (**neurotransmitter atoms and molecules**) and electrical (**charge**)

3

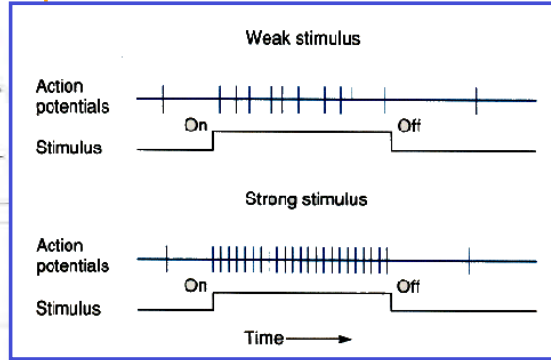
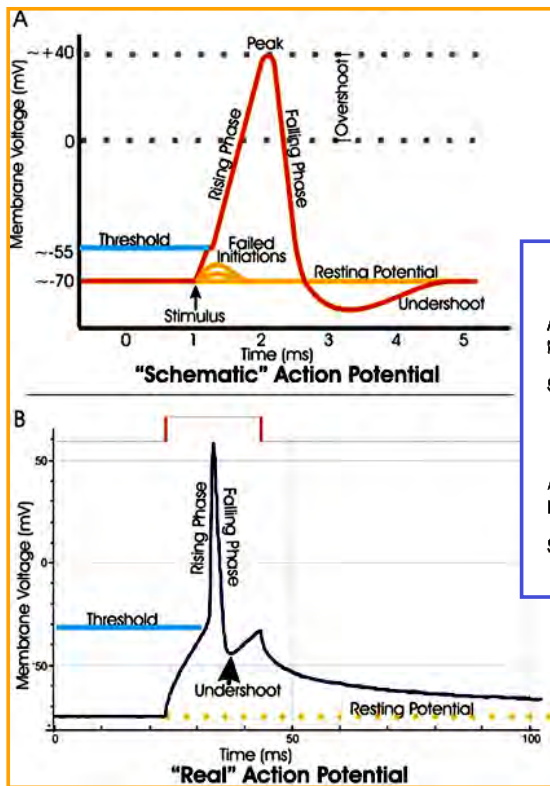
Activation Input to Soma Causes Change in Output Potential

- **Stimulus from**
 - Receptors
 - Other neurons
 - Muscle cells
 - Pacemakers (c.g. cardiac sino-atrial node)
- **When membrane potential of neuronal cell exceeds a threshold**
 - Cell is polarized
 - **Action potential** pulse is transmitted from the cell
 - Activity measured by **amplitude** and **firing frequency** of pulses
 - **Saltatory conduction** along axon
 - **Myelin Schwann cells** insulate axon
 - Signal boosted at **Nodes of Ranvier**
- **Cell depolarizes and potential returns to rest**



4

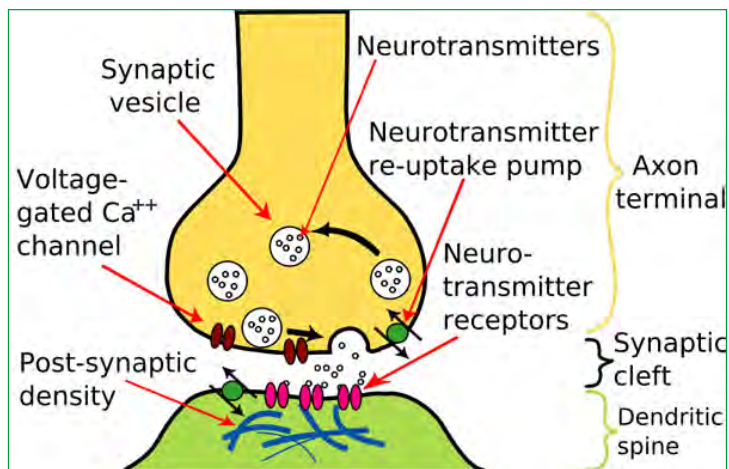
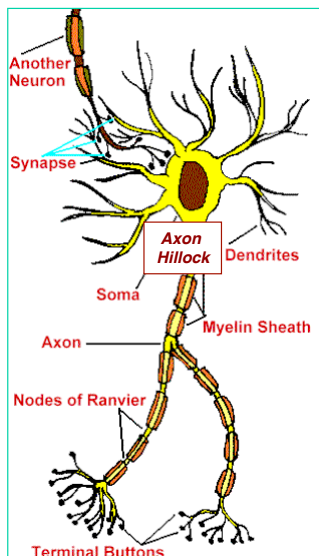
Neural Action Potential



- **Maximum Firing Rate: 500/sec**
- **Refractory Period: Minimum time increment between action potential firing ~ 1-2 msec**

5

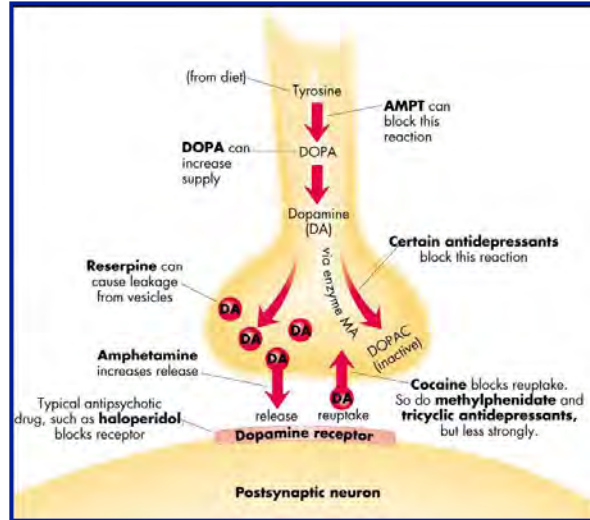
Electrochemical Signaling at Axon Hillock and Synapse



6

Synaptic Strength Can Be Increased or Decreased by Externalities

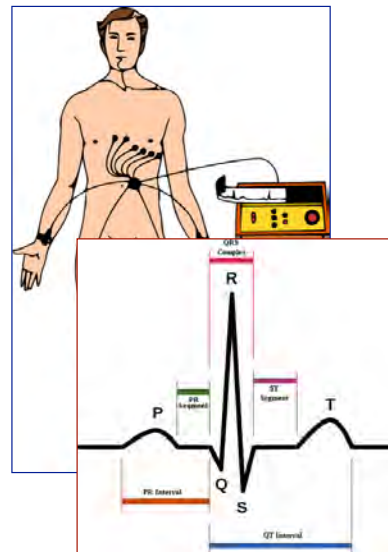
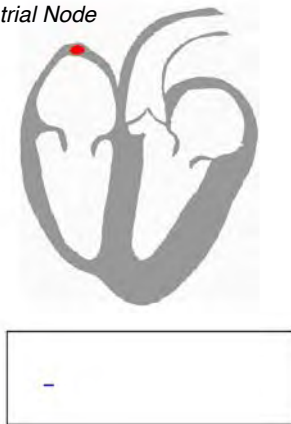
- **Synapses: learning elements of the nervous system**
 - Action potentials enhanced or inhibited
 - Chemicals can modify signal transfer
 - Potentiation of pre- and post-synaptic cells
- **Adaptation/Learning (potentiation)**
 - Short-term
 - Long-term



7

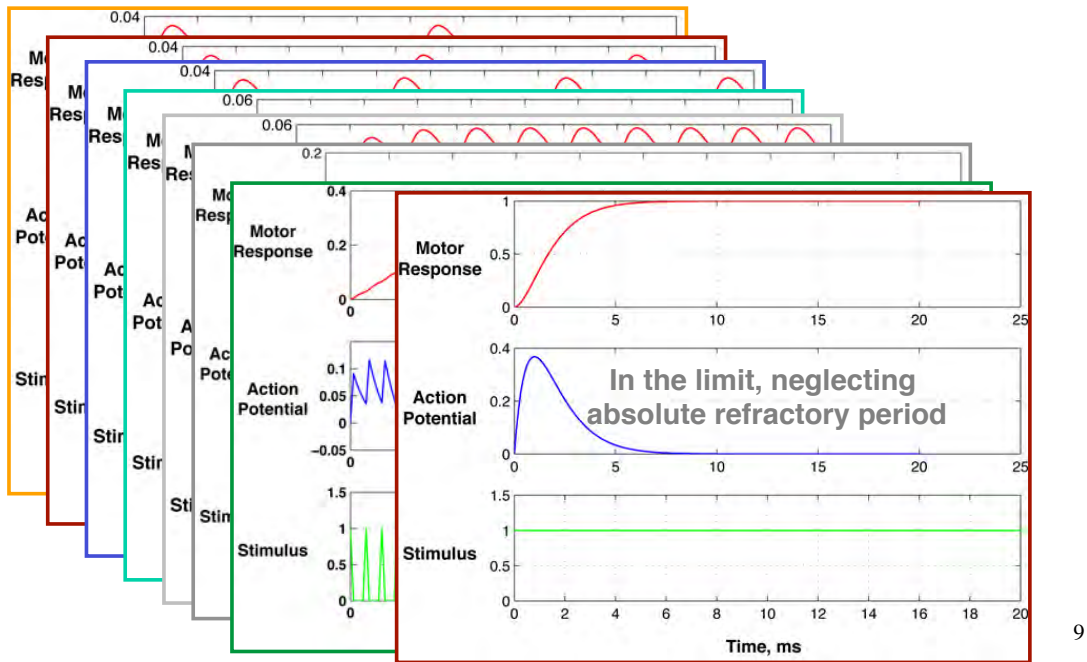
Cardiac Pacemaker and EKG Signals

Pacemaker Cell:
Sinoatrial Node

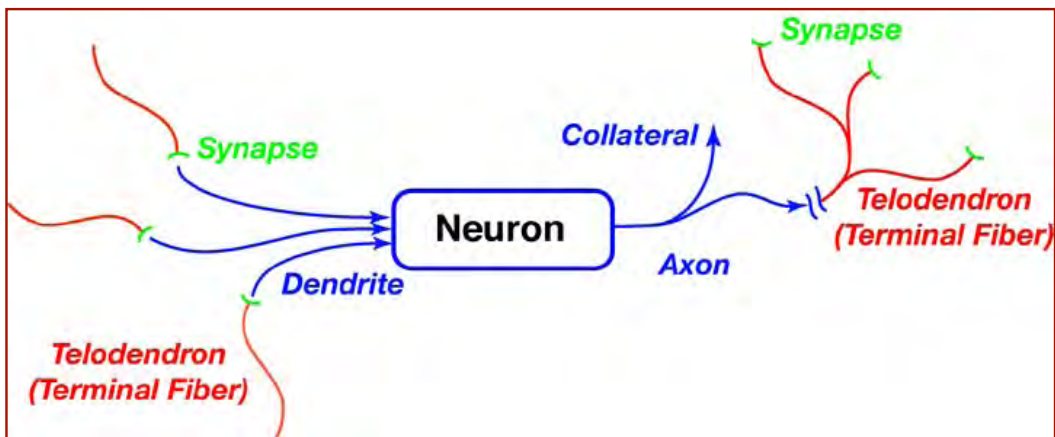


8

Impulse, Pulse-Train, and Step Response of LTI 2nd-Order Neural Model



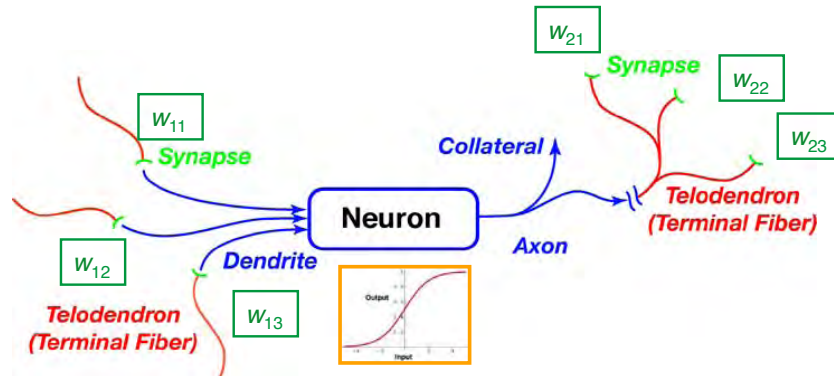
Multipolar Neuron



Mathematical Model of Neuron Components

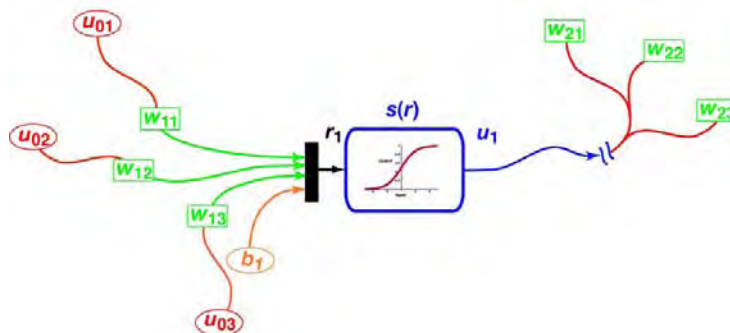
Synapse effects represented by **weights (gains or multipliers)**

Neuron firing frequency is modeled by **linear gain or nonlinear element**



11

The Neuron Function



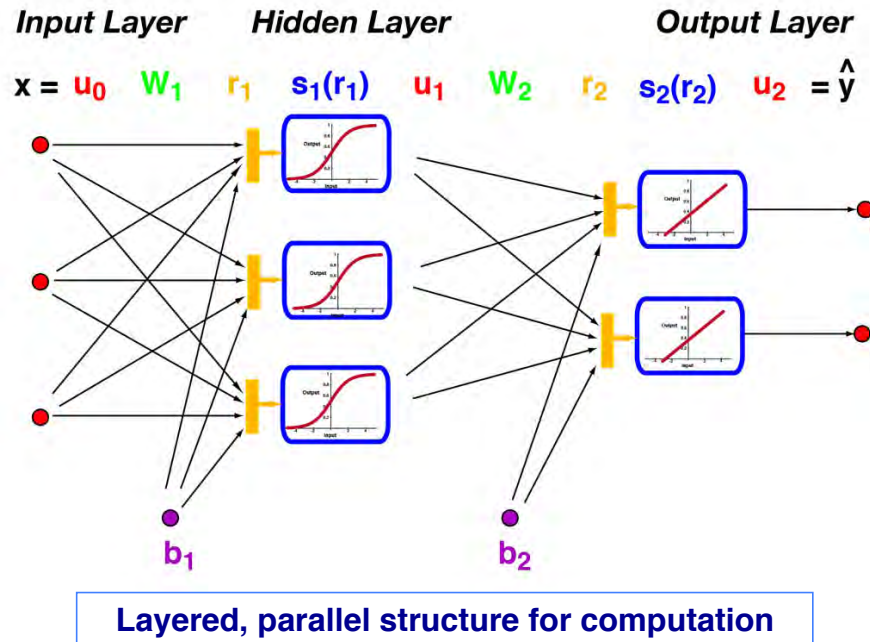
- **Vector input, \mathbf{u} , to a single neuron**
 - Sensory input or output from upstream neurons
- **Linear operation produces scalar, r**
- **Add bias, b , for zero adjustment**
- **Scalar output, u , of a single neuron (or node)**
 - Scalar linear or nonlinear operation, $s(r)$

$$r = \mathbf{w}^T \mathbf{u} + b$$

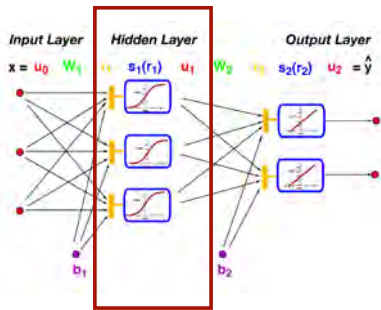
$$u = s(r)$$

12

“Shallow” Neural Network



13



Input-Output Characteristics of a Neural Network Layer

- **Single hidden layer**
 - **Number of inputs = n**
 - $\dim(u) = (n \times 1)$
 - **Number of nodes = m**
 - $\dim(r) = \dim(b) = \dim(s) = (m \times 1)$

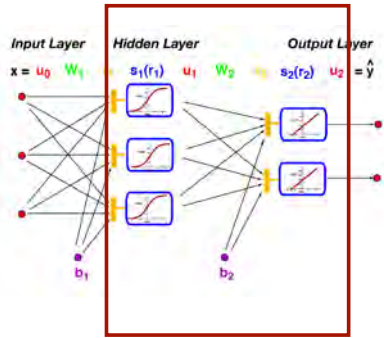
$$\mathbf{r} = \mathbf{W}\mathbf{u} + \mathbf{b}$$

$$\mathbf{u} = \mathbf{s}(\mathbf{r})$$

$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \cdots \\ \mathbf{w}_n^T \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1n} \\ \cdots & \cdots & \cdots & \cdots \\ w_{m1} & w_{m2} & \cdots & w_{mn} \end{bmatrix}$$

14

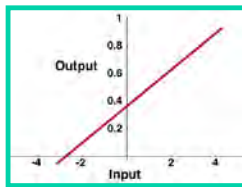
Two-Layer Network



- **Two layers**
 - **Node functions may be different, e.g.,**
 - Sigmoid hidden layer
 - Linear output layer
 - **Number of nodes in each layer need not be the same**
- **Input sometimes labeled as layer**

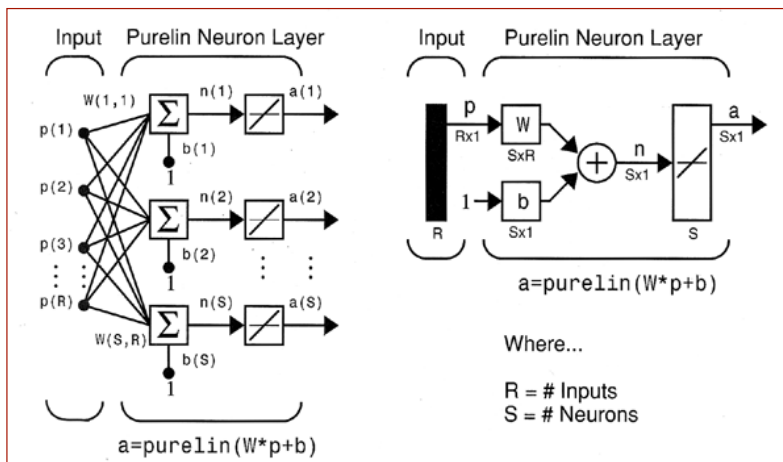
$$\begin{aligned}
 \mathbf{y} &= \mathbf{u}_2 \\
 &= \mathbf{s}_2(\mathbf{r}_2) = \mathbf{s}_2(\mathbf{W}_2 \mathbf{u}_1 + \mathbf{b}_2) \\
 &= \mathbf{s}_2[\mathbf{W}_2 \mathbf{s}_1(\mathbf{r}_1) + \mathbf{b}_2] \\
 &= \mathbf{s}_2[\mathbf{W}_2 \mathbf{s}_1(\mathbf{W}_1 \mathbf{u}_0 + \mathbf{b}_1) + \mathbf{b}_2] \\
 &= \mathbf{s}_2[\mathbf{W}_2 \mathbf{s}_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2]
 \end{aligned}$$

15



Linear Neural Network

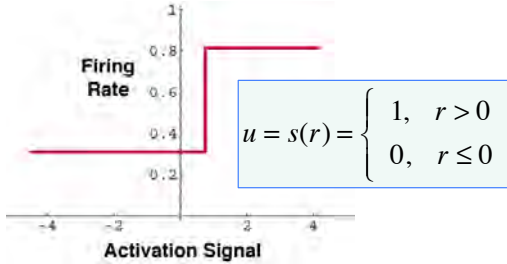
- **Outputs provide linear scaling of inputs**
- **Equivalent to matrix transformation of a vector, $\mathbf{y} = \mathbf{W}\mathbf{x} + \mathbf{b}$**
- **Easy to train (left pseudoinverse, TBD)**
- **MATLAB symbology**



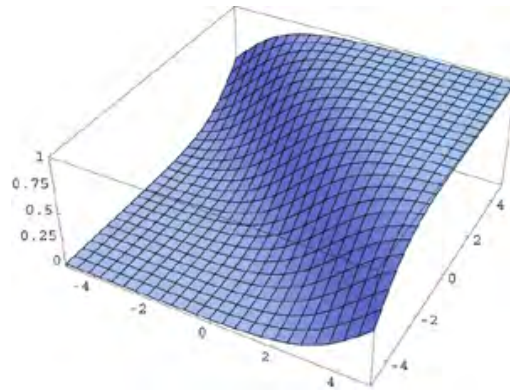
16

Idealizations of Nonlinear Neuron Input-Output Characteristic

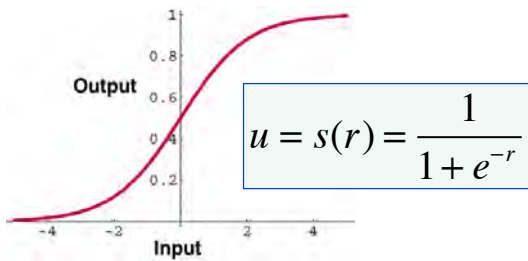
Step function ("Perceptron")



Sigmoid with two inputs, one output



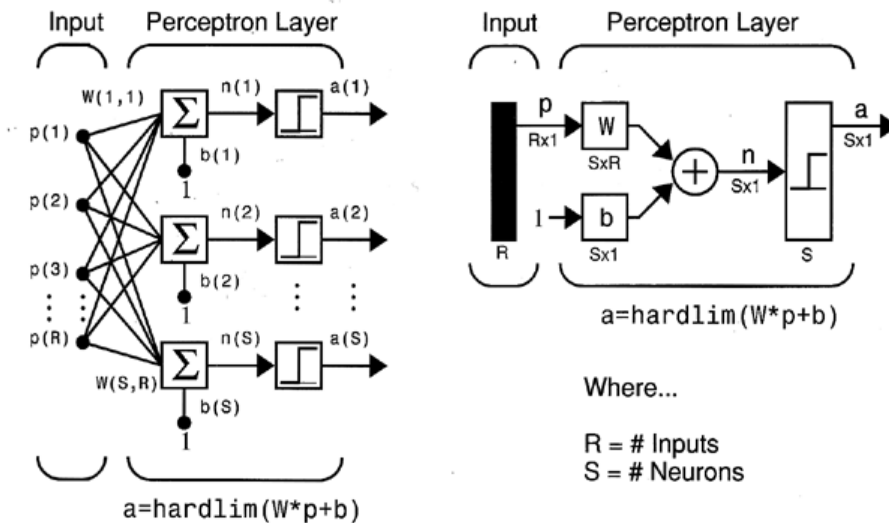
Logistic sigmoid function



$$u = s(r) = \frac{1}{1 + e^{-(w_1 r_1 + w_2 r_2 + b)}}$$

17

Perceptron Neural Network

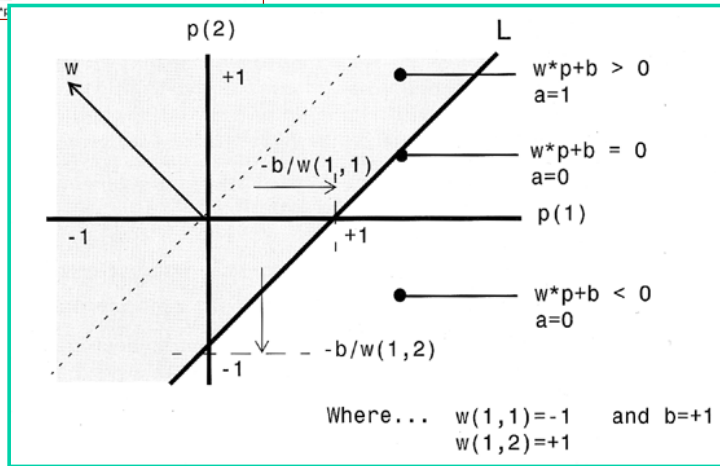
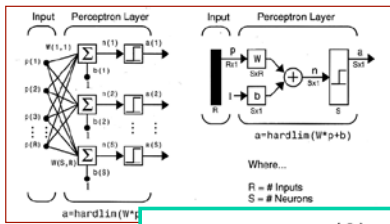


Each node is a step function

Weighted sum of features is fed to each node

Each node produces a linear classification of the input space

Perceptron Neural Network



Weights adjust slopes
Biases adjust zero crossing points

19

Single-Layer, Single-Node Perceptron Discriminants

Perceptron Function

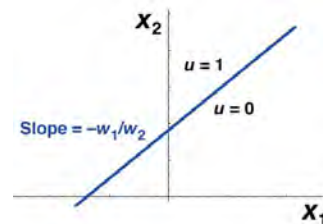
$$u = s(\mathbf{w}^T \mathbf{x} + b) = \begin{cases} 1, & (\mathbf{w}^T \mathbf{x} + b) > 0 \\ 0, & (\mathbf{w}^T \mathbf{x} + b) \leq 0 \end{cases}$$

Two inputs, single step function Discriminant

$$0 = w_1 x_1 + w_2 x_2 + b$$

$$x_2 = \frac{-1}{w_2} (w_1 x_1 + b)$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

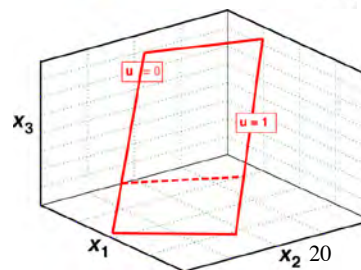


Three inputs, single step function Discriminant

$$0 = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$

$$x_3 = \frac{-1}{w_3} (w_1 x_1 + w_2 x_2 + b)$$

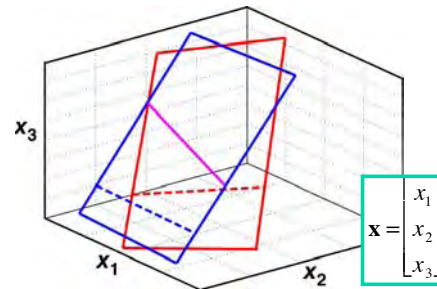
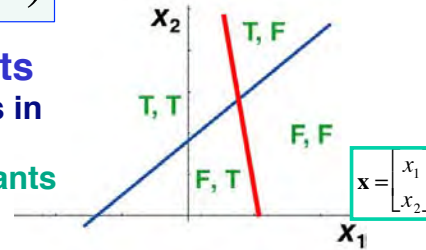
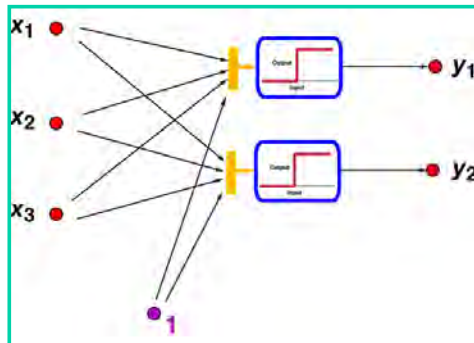
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$



Single-Layer, Multi-Node Perceptron Discriminants

$$\mathbf{u} = s(\mathbf{W}\mathbf{x} + \mathbf{b})$$

- Multiple inputs, nodes, and outputs
 - More inputs lead to more dimensions in discriminants
 - More outputs lead to more discriminants




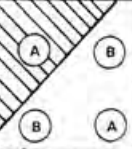


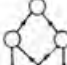
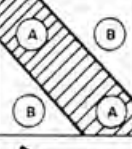



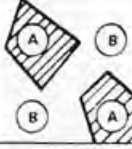


21

Multi-Layer Perceptrons Can Classify With Boundaries or Clusters

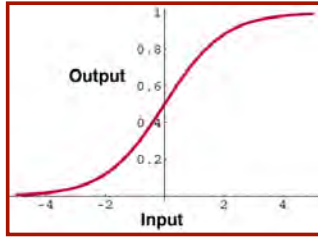
Classification capability of multi-layer perceptrons

Classifications of classifications

Open or closed regions

STRUCTURE	TYPES OF DECISION REGIONS	EXCLUSIVE OR PROBLEM	CLASSES WITH MESHED REGIONS	MOST GENERAL REGION SHAPES
 SINGLE-LAYER	HALF PLANE BOUNDED BY HYPERPLANE			
 TWO-LAYER	CONVEX OPEN OR CLOSED REGIONS			
 THREE-LAYER	ARBITRARY (Complexity Limited By Number of Nodes)			

22



Sigmoid Activation Functions

- Alternative sigmoid functions

- Logistic function: 0 to 1
- Hyperbolic tangent: -1 to 1
- Augmented ratio of squares: 0 to 1

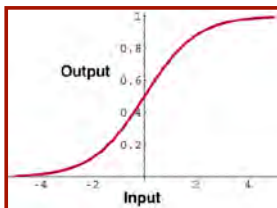
$$u = s(r) = \frac{1}{1 + e^{-r}}$$

$$u = s(r) = \tanh r = \frac{1 - e^{-2r}}{1 + e^{-2r}}$$

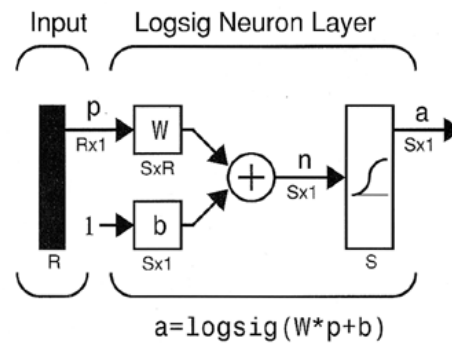
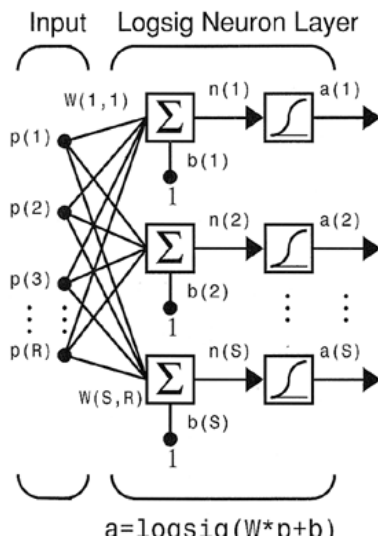
- Smooth nonlinear functions that limit extreme values in output

$$u = s(r) = \frac{r^2}{1 + r^2}$$

23



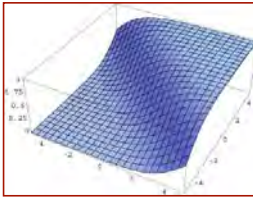
Single-Layer Sigmoid Neural Network



Where...

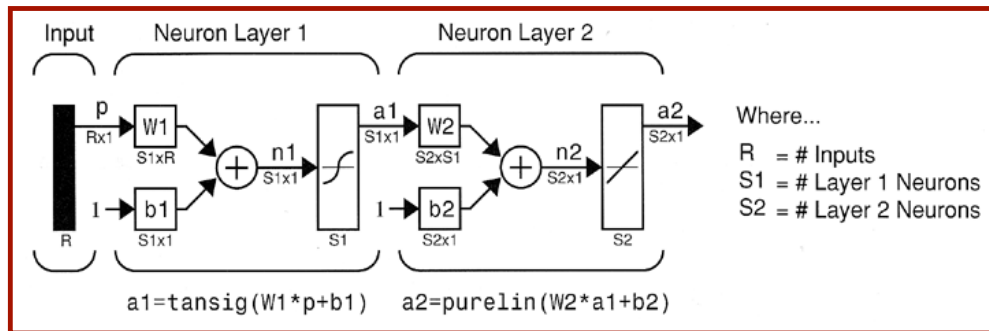
R = # Inputs
S = # Neurons

24

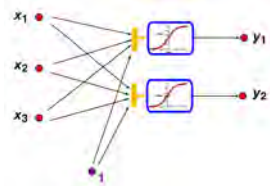


Fully Connected Two-Layer (Single-Hidden-Layer) Sigmoid Layer

- **Sufficient to approximate any continuous function**
- **All nodes of one layer connected to all nodes of adjacent layers**
- **Typical sigmoid network contains**
 - **Single sigmoid hidden layer (nonlinear fit)**
 - **Single linear output layer (scaling)**



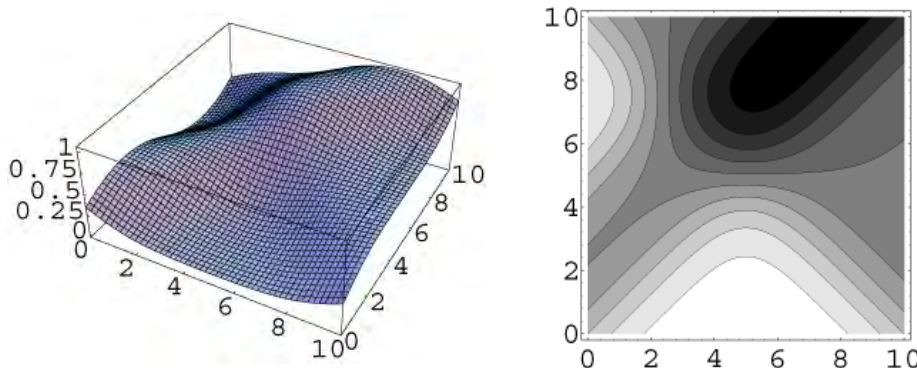
25



Typical Output for Two-Sigmoid Network

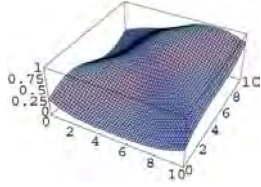
Classification is not limited to linear discriminants

2 Inputs, Single Output



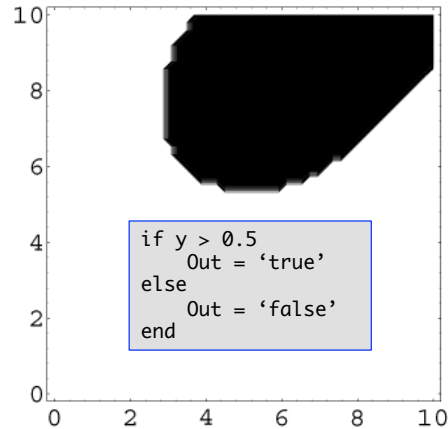
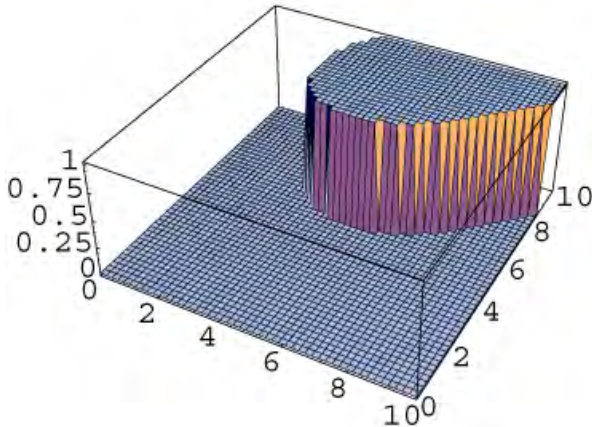
Sigmoid network can approximate a continuous nonlinear function to arbitrary accuracy with a *single hidden layer*

26



Thresholded Neural Network Output

Threshold gives “yes/no” output

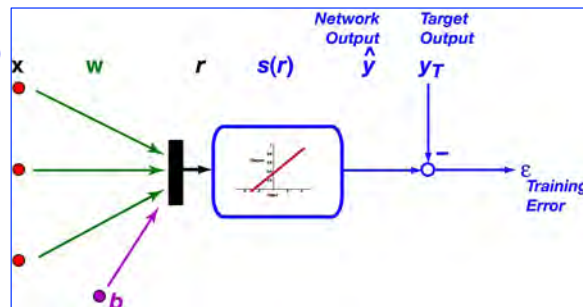


27

Least-Squares Training Example: Single Linear Neuron

- Training set (n members)
 - Target outputs, \mathbf{y}_T ($1 \times n$)
 - m Features (inputs), \mathbf{X} ($m \times n$)

$$\begin{bmatrix} \mathbf{y}_T \\ \mathbf{X} \end{bmatrix} = \begin{bmatrix} y_{T_1} & y_{T_2} & \dots & y_{T_n} \\ \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{bmatrix}_1 & \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{bmatrix}_2 & \dots & \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{bmatrix}_n \end{bmatrix}$$



- Network output, single input

$$\hat{y}_j = r_j = \hat{\mathbf{w}}^T \mathbf{x}_j + \hat{b}$$

- Quadratic error cost

- Training error

$$\epsilon_j = \hat{y}_j - y_T$$

$$J = \frac{1}{2} \sum_{j=1}^n \epsilon_j^2 = \frac{1}{2} \sum_{j=1}^n (\hat{y}_j - y_T)^2 = \frac{1}{2} \sum_{j=1}^n (\hat{y}_j^2 - 2\hat{y}_j y_T + y_T^2)$$

Note: This is an introduction to least-squares **back-propagation training**. Training of a linear neuron more readily accomplished using left pseudoinverse (Lec. 21).

28

Linear Neuron Gradient

$$\hat{y}_j = r_j = \mathbf{w}^T \mathbf{x}_j + b$$

$$\frac{d\hat{y}_j}{dr_j} = 1$$

$$\varepsilon_j = \hat{y}_j - y_T$$

$$J = \frac{1}{2} \sum_{j=1}^n \varepsilon_j^2 = \frac{1}{2} \sum_{j=1}^n (\hat{y}_j - y_T)^2 = \frac{1}{2} \sum_{j=1}^n (\hat{y}_j^2 - 2\hat{y}_j y_T + y_T^2)$$

- **Training (control) parameter, \mathbf{p}**
 - Input weights, \mathbf{w} ($n \times 1$)
 - Bias, b (1×1)

$$\mathbf{p} = \begin{bmatrix} p_1 \\ p_2 \\ \dots \\ p_{n+1} \end{bmatrix} \triangleq \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$$

- **Optimality condition**

$$\frac{\partial J}{\partial \mathbf{p}} = \mathbf{0}$$

- **Gradient**

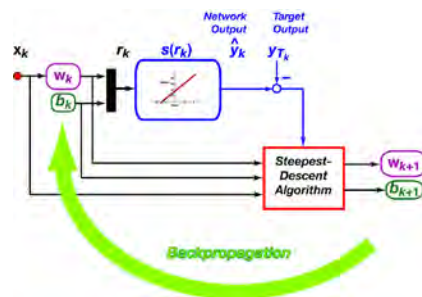
$$\frac{\partial J}{\partial \mathbf{p}} = \frac{1}{2} \sum_{j=1}^n (\hat{y}_j - y_T) \frac{\partial y_j}{\partial \mathbf{p}} = \frac{1}{2} \sum_{j=1}^n (\hat{y}_j - y_T) \frac{\partial y_j}{\partial r_j} \frac{\partial r_j}{\partial \mathbf{p}}$$

where

$$\frac{\partial r_j}{\partial \mathbf{p}} = \begin{bmatrix} \frac{\partial r_j}{\partial p_1} & \frac{\partial r_j}{\partial p_2} & \dots & \frac{\partial r_j}{\partial p_{n+1}} \end{bmatrix} = \frac{\partial (\mathbf{w}^T \mathbf{x}_j + b)}{\partial \mathbf{p}} = \begin{bmatrix} \mathbf{x}_j^T & 1 \end{bmatrix}$$

29

Steepest-Descent (Back-propagation) Learning for a Single Linear Neuron



Gradient

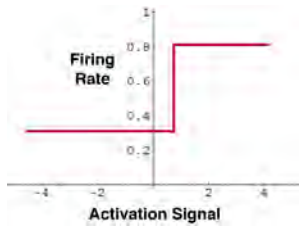
$$\frac{\partial J}{\partial \mathbf{p}} = \frac{1}{2} \sum_{j=1}^n (\hat{y}_j - y_T) \begin{bmatrix} \mathbf{x}_j^T & 1 \end{bmatrix} = \frac{1}{2} \sum_{j=1}^n [(\mathbf{w}^T \mathbf{x}_j + b) - y_T] \begin{bmatrix} \mathbf{x}_j^T & 1 \end{bmatrix}$$

Steepest-descent algorithm

$$\mathbf{p}_{k+1} = \mathbf{p}_k - \eta \left(\frac{\partial J}{\partial \mathbf{p}} \right)_k^T$$

η = learning rate
 k = iteration index(epoch)

30



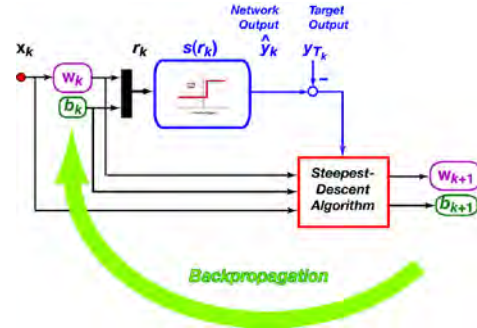
Steepest-Descent Algorithm for a Single-Step Perceptron

Neuron output is discontinuous

$$\hat{y} = s(r) = \begin{cases} 1, & r > 0 \\ 0, & r \leq 0 \end{cases}$$

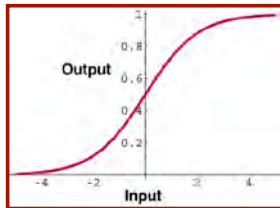
Binary target output
 $y_T = 0$ or 1 , for classification

$$(\hat{y}_{jk} - y_{T_k}) = \begin{cases} 1, & \hat{y}_{jk} = 1, y_{T_k} = 0 \\ 0, & \hat{y}_{jk} = y_{T_k} \\ -1, & \hat{y}_{jk} = 0, y_{T_k} = 1 \end{cases}$$



$$\begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}_{k+1} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}_k - \eta \sum_{j=1}^N [\hat{y}_{jk} - y_{T_k}] \begin{bmatrix} \mathbf{x}_j \\ 1 \end{bmatrix}$$

31



Training Variables for a Single Sigmoid Neuron

Neuron output is continuous

$$\begin{aligned} \hat{y} = s(r) &= \frac{1}{1 + e^{-r}} \\ &= s(\mathbf{w}^T \mathbf{x} + b) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}} \end{aligned}$$

Training error and quadratic error cost

$$\begin{aligned} \varepsilon_j &= \hat{y}_j - y_T \\ J &= \frac{1}{2} \sum_{j=1}^n \varepsilon_j^2 = \frac{1}{2} \sum_{j=1}^n (\hat{y}_j - y_T)^2 = \frac{1}{2} \sum_{j=1}^n (\hat{y}_j^2 - 2\hat{y}_j y_T + y_T^2) \end{aligned}$$

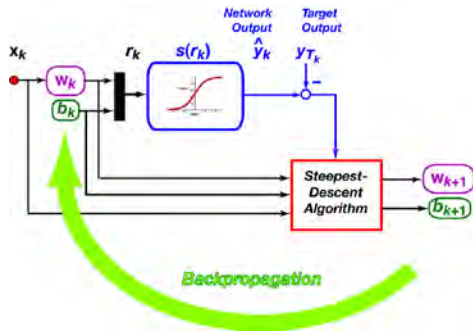
Neuron output sensitivity to input

$$\begin{aligned} \frac{d\hat{y}}{dr} &= \frac{ds(r)}{dr} = \frac{e^{-r}}{(1 + e^{-r})^2} = e^{-r} s^2(r) \\ &= [(1 + e^{-r}) - 1] s^2(r) = \left[\frac{1 - s(r)}{s(r)} \right] s^2(r) \end{aligned}$$

$$\frac{d\hat{y}}{dr} = [1 - s(r)] s(r) = (1 - \hat{y}) \hat{y}$$

32

Back-Propagation Training of a Single Sigmoid Neuron



$$\frac{\partial J}{\partial \mathbf{p}} = \frac{1}{2} \sum_{j=1}^N (\hat{y}_j - y_{Tj}) \frac{\partial \hat{y}_j}{\partial r} \frac{\partial r}{\partial \mathbf{p}}$$

$$\mathbf{p}_{k+1} = \mathbf{p}_k - \eta \left(\frac{\partial J}{\partial \mathbf{p}} \right)_k$$

or

where

$$r = \mathbf{w}^T \mathbf{x} + b$$

$$\frac{d\hat{y}}{dr} = (1 - \hat{y})\hat{y}$$

$$\frac{\partial r}{\partial \mathbf{p}} = \begin{bmatrix} \mathbf{x}^T & 1 \end{bmatrix}$$

$$\begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}_{k+1} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}_k - \eta \sum_{j=1}^N \left\{ [\hat{y}_{jk} - y_{Tj}] (1 - \hat{y}_k) \hat{y}_k \begin{bmatrix} \mathbf{x}_j \\ 1 \end{bmatrix} \right\}_k$$

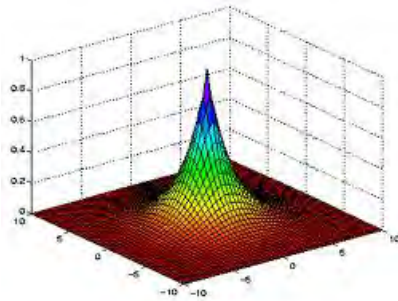
See Supplemental Material for training multiple sigmoids

Radial Basis Function

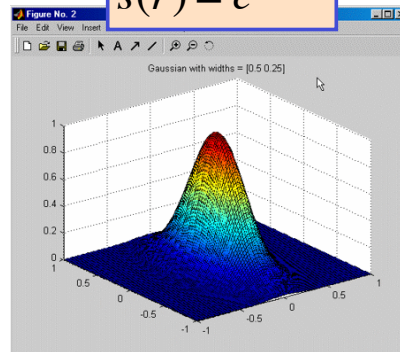
Unimodal, axially symmetric function, e.g., exponential

$$s(r) = e^{-|ar|^n}, \quad r = \sqrt{(\mathbf{x} - \mathbf{x}_{center})^T (\mathbf{x} - \mathbf{x}_{center})}$$

$$s(r) = e^{-|ar|}$$



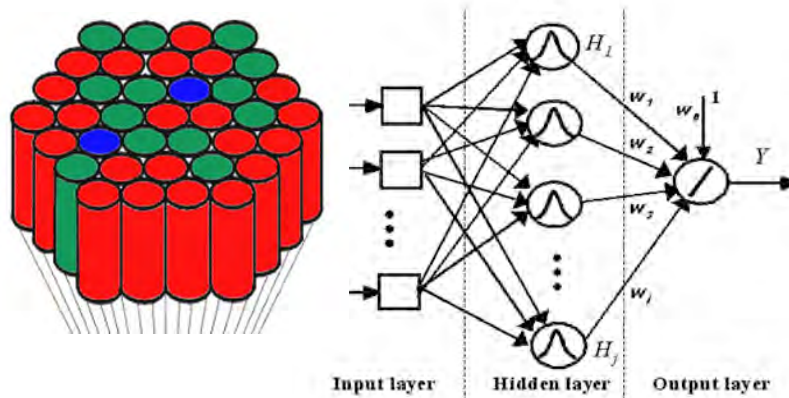
$$s(r) = e^{-(ar)^2}$$



Network mimics stimulus field of a neuron receptor, e.g., retina

Radial Basis Function Network

Array of RBFs typically centered on a fixed grid



http://en.wikipedia.org/wiki/Radial_basis_function_network

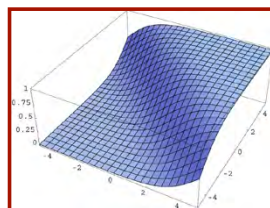
35

Sigmoid vs. Radial Basis Function Node

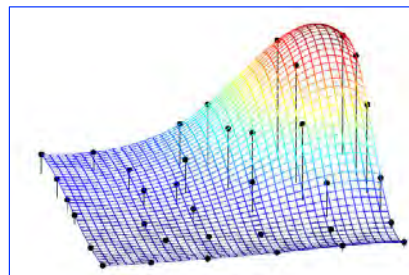
- Considerations for selecting the basis function
 - Prior knowledge of surface to be approximated
 - Global vs. compact support
 - Number of neurons required
 - Training and untraining issues

Sigmoid function

$$s(r) = \frac{1}{1 + e^{-r}}$$

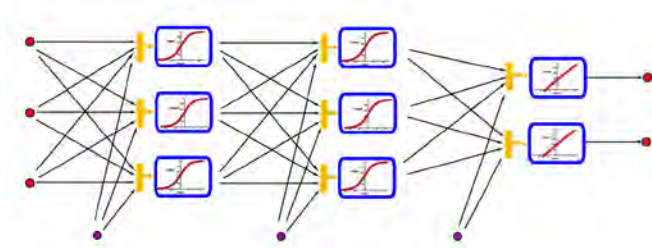


Radial basis functions



36

“Deep” Sigmoid Network



- Multiple hidden and “visible” layers can improve accuracy in image processing and language translation
- Problem of the “vanishing gradient” in training
- One solution: *Convolutional neural network* of neuron input/output by incremental training
 - Pooling or clustering signals between layers (*TBD*)
 - Limited receptive fields for filter (or kernel) nodes
 - Node is activated only when input is within pre-determined bounds (see *CMAC*, Lecture 19)

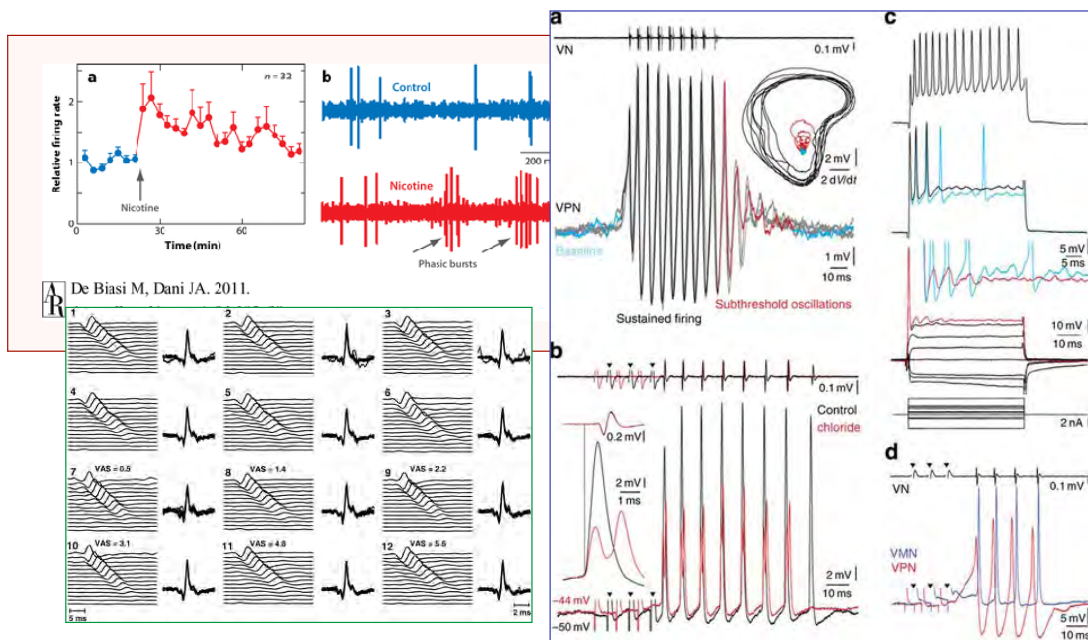
37

*Next Time:
More on Neural Networks*

Supplementary Material

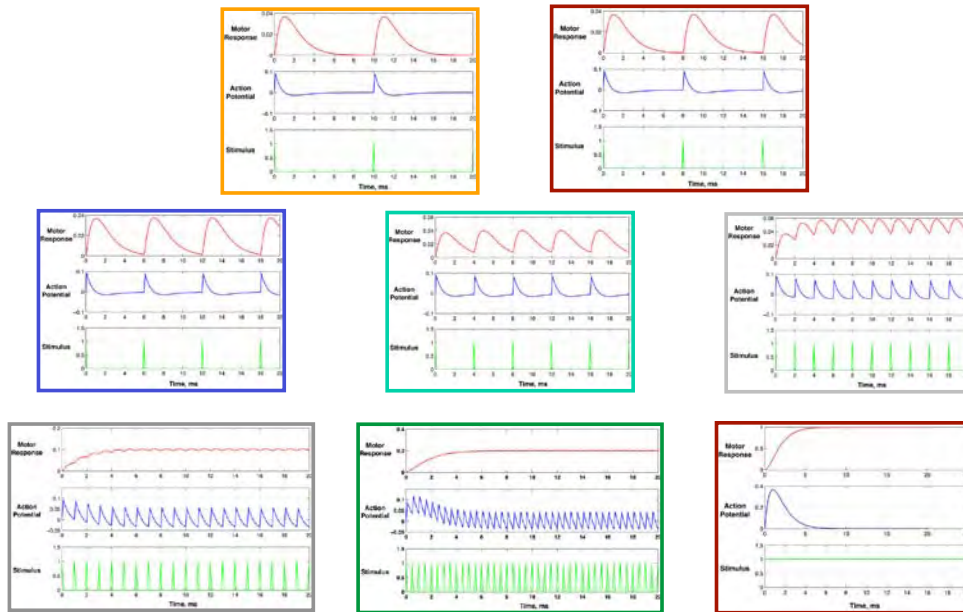
39

Some Recorded Action Potential Pulse Trains



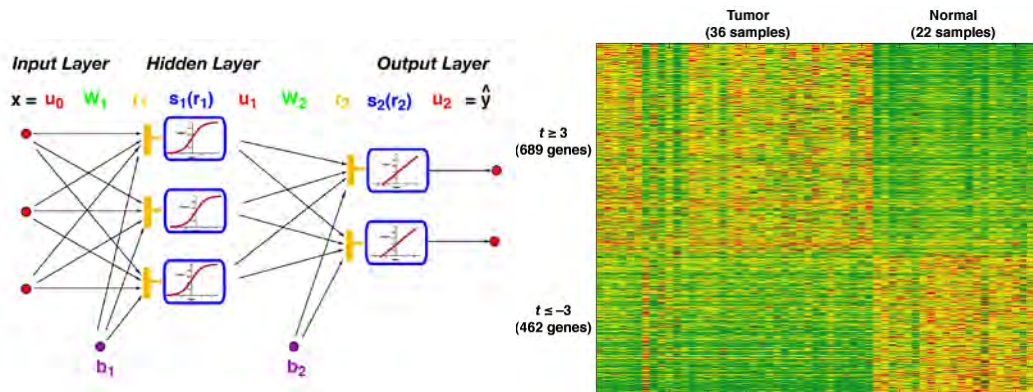
40

Impulse, Pulse-Train, and Step Response of a LTI 2nd-Order Neural Model



41

Microarray Training Set

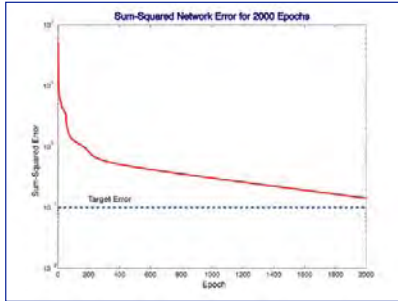


$\begin{bmatrix} \text{Identifier} \\ \mathbf{y}_T \\ \mathbf{X} \end{bmatrix} =$	<i>Sample 1</i>	<i>Sample 2</i>	<i>Sample 3</i>	...	<i>Sample n-1</i>	<i>Sample n</i>
	<i>Tumor</i>	<i>Tumor</i>	<i>Tumor</i>	...	<i>Normal</i>	<i>Normal</i>
	<i>Gene A Level</i>	<i>Gene A Level</i>	<i>Gene A Level</i>	...	<i>Gene A Level</i>	<i>Gene A Level</i>
	<i>Gene B Level</i>	<i>Gene B Level</i>	<i>Gene B Level</i>	...	<i>Gene B Level</i>	<i>Gene B Level</i>

	<i>Gene m Level</i>	<i>Gene m Level</i>	<i>Gene m Level</i>	...	<i>Gene m Level</i>	<i>Gene m Level</i>

42

Neural Network Training Results: Tumor/Normal Classification, 8 Genes, 4 Nodes



- Training begins with a random set of weights
- Adjustable parameters
 - Learning rate
 - Target error
 - Maximum # of epochs
- Non-unique sets of trained weights

Binary network output (0,1) after rounding

Zero classification errors

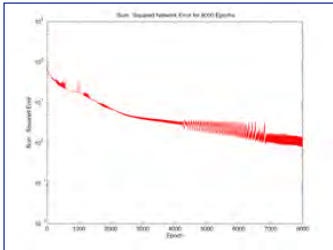
Classification =

Columns 1 through 13	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Columns 14 through 26	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Columns 27 through 39	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Columns 40 through 52	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Columns 53 through 62	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	45

Neural Network Training Results: Tumor Stage/Normal Classification 8 Genes, 16 Nodes

- Colon cancer classification
 - 0 = Normal
 - 1 = Adenoma
 - 2 = A Tumor
 - 3 = B Tumor
 - 4 = C Tumor
 - 5 = D Tumor

Scalar network output with varying magnitude



Target =

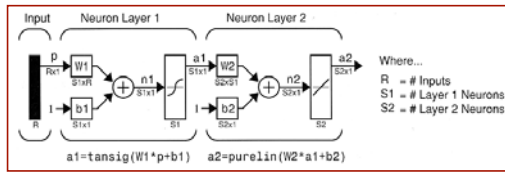
[2	1	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	4
4	4	4	4	4	4	4	4	4	5	5	5	5	5	5	5	5
5	5	5	5	5	1	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

One classification error

Classification =

Columns 1 through 13	2	1	3	3	3	3	3	3	3	3	3	3	3	3	3	3
Columns 14 through 26	3	3	3	3	3	3	3	4	4	5	4	4	4	4	4	4
Columns 27 through 39	4	4	4	5	5	5	5	5	5	5	5	5	5	1	0	0
Columns 40 through 52	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Columns 53 through 60	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Training a Sigmoid Network



Two parameter vectors for 2-layer network

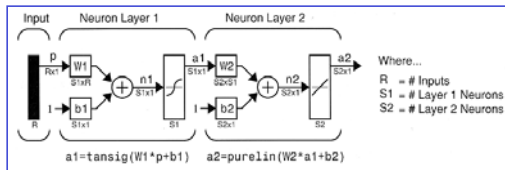
$$\mathbf{p}_{1,2} = \begin{bmatrix} \mathbf{w} \\ \mathbf{b} \end{bmatrix}_{1,2} = \begin{bmatrix} p_1 \\ p_2 \\ \dots \\ p_{n+1} \end{bmatrix}_{1,2}$$

Output vector

$$\begin{aligned} \hat{\mathbf{y}} &= \mathbf{u}_2 \\ &= \mathbf{s}_2(\mathbf{r}_2) = \mathbf{s}_2(\mathbf{W}_2 \mathbf{u}_1 + \mathbf{b}_2) \\ &= \mathbf{s}_2[\mathbf{W}_2 \mathbf{s}_1(\mathbf{r}_1) + \mathbf{b}_2] \\ &= \mathbf{s}_2[\mathbf{W}_2 \mathbf{s}_1(\mathbf{W}_1 \mathbf{u}_0 + \mathbf{b}_1) + \mathbf{b}_2] \\ &= \mathbf{s}_2[\mathbf{W}_2 \mathbf{s}_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2] \end{aligned}$$

47

Training a Sigmoid Network



$$\mathbf{p}_{1,2,k} = \mathbf{p}_{1,2,k} - \eta \left(\frac{\partial J}{\partial \mathbf{p}_{1,2}} \right)_k^T$$

where

$$\frac{\partial J}{\partial \mathbf{p}_{1,2}} = (\hat{\mathbf{y}} - \mathbf{y}_T) \frac{\partial \mathbf{y}}{\partial \mathbf{p}_{1,2}} = (\hat{\mathbf{y}} - \mathbf{y}_T) \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{r}_{1,2}} \frac{\partial \mathbf{r}_{1,2}}{\partial \mathbf{p}_{1,2}}$$

$$\begin{aligned} \mathbf{r}_{1,2} &= \mathbf{W}_{1,2} \mathbf{u}_{0,1} + \mathbf{b}_{1,2} \\ \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{r}_2} &= \mathbf{I}; \quad \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{r}_1} = \begin{bmatrix} (1-\hat{y}_1)\hat{y}_1 & 0 & \dots & 0 \\ 0 & (1-\hat{y}_2)\hat{y}_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & (1-\hat{y}_n)\hat{y}_n \end{bmatrix} \\ \frac{\partial \mathbf{r}_1}{\partial \mathbf{p}_1} &= \begin{bmatrix} \mathbf{x}^T & 1 \end{bmatrix}; \quad \frac{\partial \mathbf{r}_2}{\partial \mathbf{p}_2} = \begin{bmatrix} \mathbf{u}_1^T & 1 \end{bmatrix} \end{aligned}$$

48

MATLAB Neural Network Toolbox

Training Algorithms

Backpropagation training functions that use Jacobian derivatives

These algorithms can be faster but require more memory than gradient backpropagation. They are also not supported on GPU hardware.

- [trainlm](#) - Levenberg-Marquardt backpropagation.
- [trainbr](#) - Bayesian Regulation backpropagation.

Backpropagation training functions that use gradient derivatives

These algorithms may not be as fast as Jacobian backpropagation. They are supported on GPU hardware with the Parallel Computing Toolbox.

- [trainbfg](#) - BFGS quasi-Newton backpropagation.
- [traincgb](#) - Conjugate gradient backpropagation with Powell-Beale restarts.
- [traincgf](#) - Conjugate gradient backpropagation with Fletcher-Reeves updates.
- [traincgp](#) - Conjugate gradient backpropagation with Polak-Ribiere updates.
- [traingd](#) - Gradient descent backpropagation.
- [traingda](#) - Gradient descent with adaptive lr backpropagation.
- [traingdm](#) - Gradient descent with momentum.
- [traingdx](#) - Gradient descent w/momentum & adaptive lr backpropagation.
- [trainoss](#) - One step secant backpropagation.
- [trainrp](#) - RPROP backpropagation.
- [trainscg](#) - Scaled conjugate gradient backpropagation.

Supervised weight/bias training functions

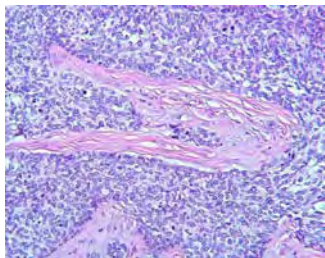
- [trainb](#) - Batch training with weight & bias learning rules.
- [trainc](#) - Cyclical order weight/bias training.
- [trainr](#) - Random order weight/bias training.
- [trains](#) - Sequential order weight/bias training.

Unsupervised weight/bias training functions

- [trainbu](#) - Unsupervised batch training with weight & bias learning rules.
- [trainru](#) - Unsupervised random order weight/bias training.

49

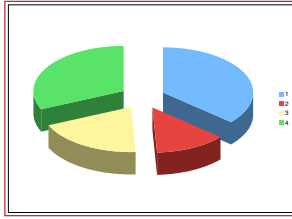
Small, Round Blue-Cell Tumor Classification Example



Desmoplastic small, round blue-cell tumors

- Childhood cancers, including
 - Ewing's sarcoma (EWS)
 - Burkitt's Lymphoma (BL)
 - Neuroblastoma (NB)
 - Rhabdomyosarcoma (RMS)
- cDNA microarray analysis presented by J. Khan, *et al.*, *Nature Medicine*, 2001, 673-679.
 - 96 transcripts chosen from 2,308 probes for training
 - 63 EWS, BL, NB, and RMS training samples
- Source of data for my analysis

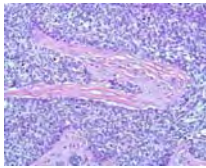
50



Overview of Present SRBCT Analysis

- **Transcript selection by *t* test**
 - 96 transcripts, 12 highest and lowest *t* values for each class
 - **Overlap with Khan set: 32 transcripts**
- **Ensemble averaging** of genes with highest and lowest *t* values in each class
- **Cross-plot** of ensemble averages
- **Classification by sigmoidal neural network**
- **Validation** of neural network
 - **Novel set simulation**
 - **Leave-one-out simulation**

51



Ranking by EWS *t* Values (Top and Bottom 12)

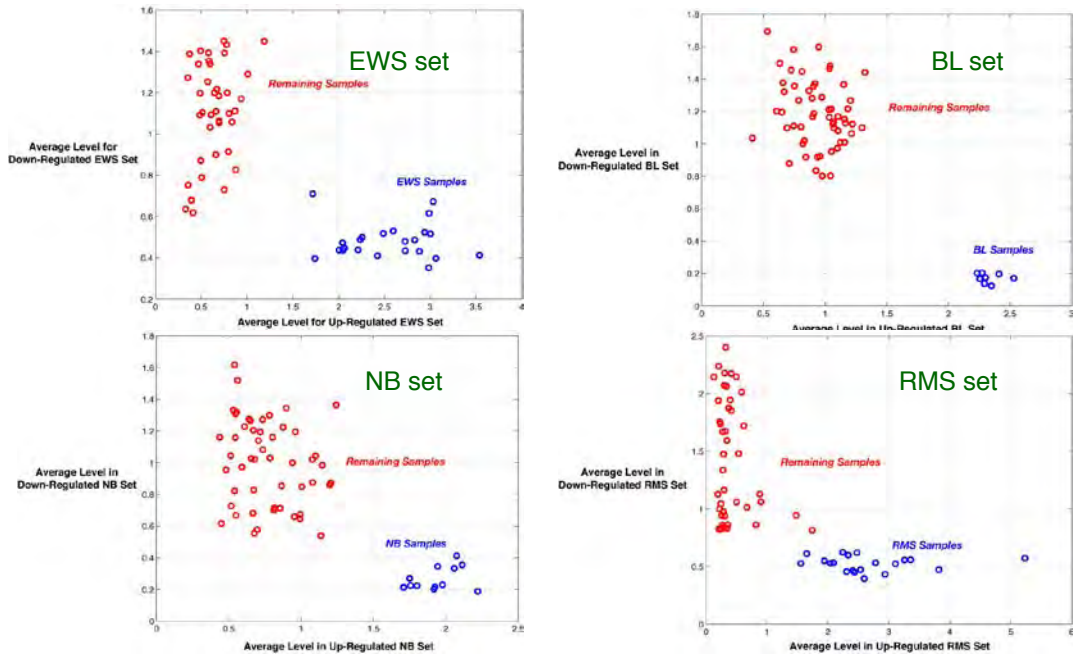
- 24 transcripts selected from 12 highest and lowest *t* values for EWS vs. remainder

Image ID	Sort by EWS <i>t</i> Value Transcript Description	EWS <i>t</i> Value	BL <i>t</i> Value	NB <i>t</i> Value	RMS <i>t</i> Value
770394	Fc fragment of IgG, receptor, transporter, alpha	12.04	-6.67	-6.17	-4.79
1435862	antigen identified by monoclonal antibodies 12E7, F21 and O13	9.09	-6.75	-5.01	-4.03
377461	caveolin 1, caveolae protein, 22kD	8.82	-5.97	-4.91	-4.78
814260	follicular lymphoma variant translocation 1	8.17	-4.31	-4.70	-5.48
491565	Chp/p300-interacting transactivator, with Glu/Asp-rich carboxy-terminal domain	7.60	-5.82	-2.62	-3.68
841641	cyclin D1 (PRAD1; parathyroid adenomatosis 1)	6.84	-9.93	0.56	-4.30
1471841	ATPase, Na ⁺ /K ⁺ transporting, alpha 1 polypeptide	6.65	-3.56	-2.72	-4.69
866702	protein tyrosine phosphatase, non-receptor type 13	6.54	-4.99	-4.07	-4.84
713922	glutathione S-transferase M1	6.17	-5.61	-5.16	-1.97
308497	KIAA0467 protein	5.99	-6.69	-6.63	-1.11
770868	NGFI-A binding protein 2 (ERG1 binding protein 2)	5.93	-6.74	-3.88	-1.21
345232	lymphotoxin alpha (TNF superfamily, member 1)	5.61	-8.05	-2.49	-1.19
786084	chromobox homolog 1 (Drosophila HP1 beta)	-5.04	-1.05	9.65	-0.62
796258	sarcoglycan, alpha (50kD dystrophin-associated glycoprotein)	-5.04	-3.31	-3.86	6.83
431397		-5.04	2.64	2.19	0.64
825411	N-acetylglucosamine receptor 1 (thyroid)	-5.06	-1.45	5.79	0.76
859359	quinone oxidoreductase homolog	-5.23	-7.27	0.78	5.40
75254	cysteine and glycine-rich protein 2 (LIM domain only, smooth muscle)	-5.30	-4.11	2.20	3.68
448386		-5.38	-0.42	3.76	0.14
68950	cyclin E1	-5.80	0.03	-1.58	5.10
774502	protein tyrosine phosphatase, non-receptor type 12	-5.80	-5.56	3.76	3.66
842820	inducible poly(A)-binding protein	-6.14	0.60	0.66	3.80
214572	ESTs	-6.39	-0.08	-0.22	4.56
295985	ESTs	-9.26	-0.13	3.24	2.95

Repeated for BL vs. remainder, NB vs. remainder, and RMS vs. remainder

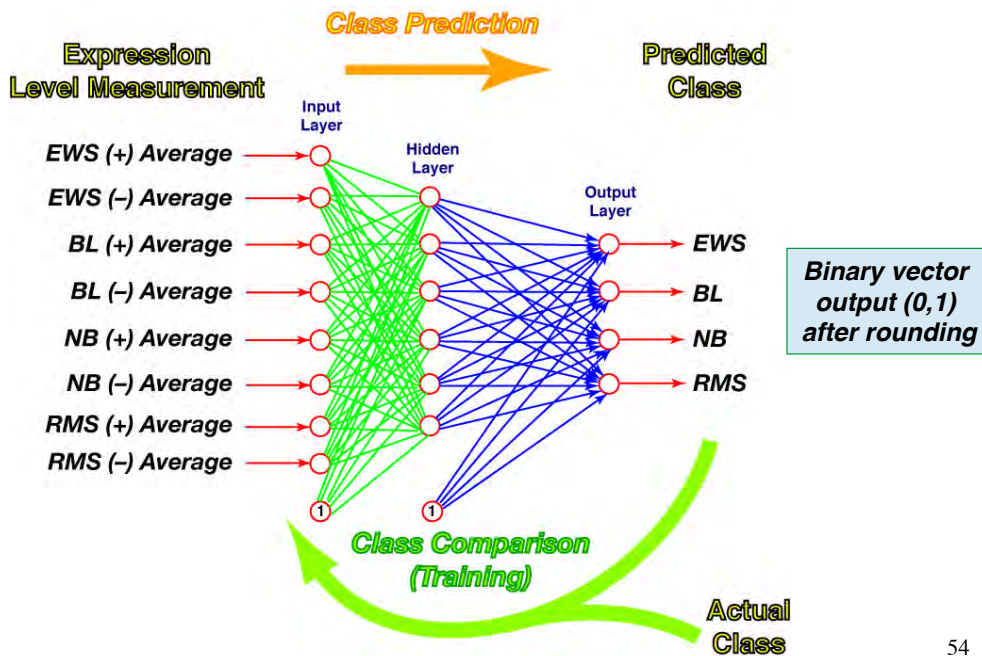
52

Clustering of SRBCT Ensemble Averages



53

SRBCT Neural Network



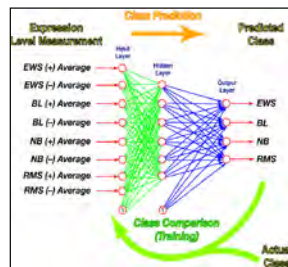
54

Neural Network Training Set

Each input row is an **ensemble average** for a transcript set, normalized in (-1,+1)

Identifier	Sample 1	Sample 2	Sample 3	...	Sample 62	Sample 63
Target Output	EWS	EWS	EWS	...	RMS	RMS
Transcript Training Set	EWS(+) <i>Average</i>	EWS(+) <i>Average</i>	EWS(+) <i>Average</i>	...	EWS(+) <i>Average</i>	EWS(+) <i>Average</i>
	EWS(-) <i>Average</i>	EWS(-) <i>Average</i>	EWS(-) <i>Average</i>	...	EWS(-) <i>Average</i>	EWS(-) <i>Average</i>
	BL(+) <i>Average</i>	BL(+) <i>Average</i>	BL(+) <i>Average</i>	...	BL(+) <i>Average</i>	BL(+) <i>Average</i>
	BL(-) <i>Average</i>	BL(-) <i>Average</i>	BL(-) <i>Average</i>	...	BL(-) <i>Average</i>	BL(-) <i>Average</i>
	NB(+) <i>Average</i>	NB(+) <i>Average</i>	NB(+) <i>Average</i>	...	NB(+) <i>Average</i>	NB(+) <i>Average</i>
	NB(-) <i>Average</i>	NB(-) <i>Average</i>	NB(-) <i>Average</i>	...	NB(-) <i>Average</i>	NB(-) <i>Average</i>
	RMS(+) <i>Average</i>	RMS(+) <i>Average</i>	RMS(+) <i>Average</i>	...	RMS(+) <i>Average</i>	RMS(+) <i>Average</i>
	RMS(-) <i>Average</i>	RMS(-) <i>Average</i>	RMS(-) <i>Average</i>	...	RMS(-) <i>Average</i>	RMS(-) <i>Average</i>

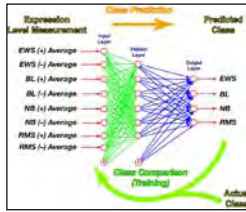
55



SRBCT Neural Network Training

- **Neural network**
 - 8 ensemble-average inputs
 - various # of sigmoidal neurons
 - 4 linear output neurons
 - 4 outputs
- **Training accuracy**
 - Train on all 63 samples
 - Test on all 63 samples
- **100% accuracy**

56

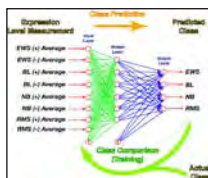


Leave-One-Out Validation of SRBCT Neural Network

- Remove a single sample
- Train on remaining samples (125 times)
- Evaluate class of the removed sample
- Repeat for each of 63 samples

- **6 sigmoids:** 99.96% accuracy (3 errors in 7,875 trials)
- **12 sigmoids:** 99.99% accuracy (1 error in 7,875 trials)

57



Novel-Set Validation of SRBCT Neural Network

- Network always chooses one of four classes (i.e., “unknown” is not an option)
- Test on 25 novel samples (400 times each)
 - 5 EWS
 - 5 BL
 - 5 NB
 - 5 RMS
 - 5 samples of unknown class

- **99.96% accuracy on first 20 novel samples** (3 errors in 8,000 trials)
- **0% accuracy on unknown classes**

58

Observations of SRBCT Classification using Ensemble Averages

- ***t* test** identified strong features for classification in this data set
 - Neural networks easily classified the four data types
 - **Caveat:** Small, round blue-cell tumors occur in different tissue types
 - Ewing' s sarcoma: **Bone tissue**
 - Burkitt' s Lymphoma: **Lymph nodes**
 - Neuroblastoma: **Nerve tissue**
 - Rhabdomyosarcoma: **Soft tissue**
- Gene expression (i.e., mRNA) level is linked to tissue difference as well as tumor genetics**

59

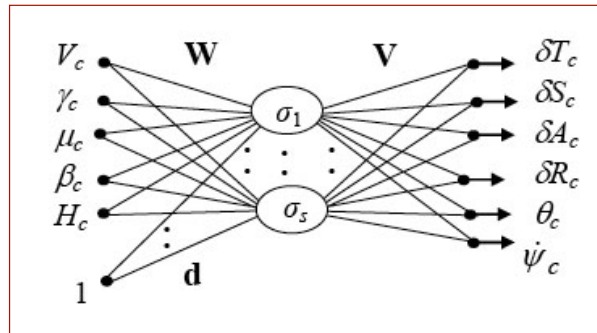
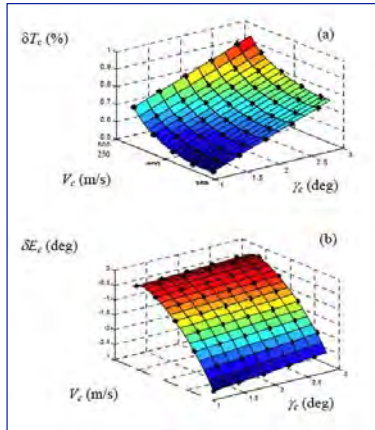
Algebraic Training of a Neural Network

Ferrari, S. and Stengel, R.,
[Smooth Function Approximation Using Neural
Networks \(pdf\)](#), *IEEE Trans. Neural Networks*, Vol. 16,
No. 1, Jan 2005, pp. 24-38 (with S. Ferrari).

60

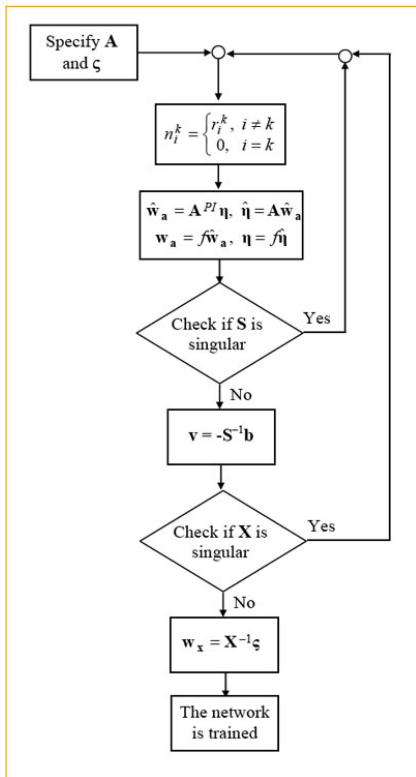
Algebraic Training for Exact Fit to a Smooth Function

- Smooth functions define equilibrium control settings at many operating points
- Neural network required to fit the functions

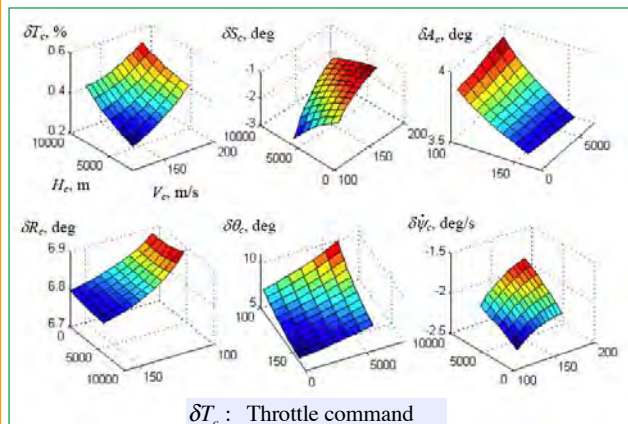


Ferrari and Stengel

61



Algorithm for Network Training



- δT_c : Throttle command
- δS_c : Spoiler command
- δA_c : Aileron command
- δR_c : Rudder command
- $\delta \theta_c$: Pitch angle command
- $\delta \dot{\psi}_c$: Yaw rate command

62

Results for Network Training

45-node example
Algorithm is considerably faster
than search methods

Algorithm:	Time (Scaled):	Flops:	Lines of code (MATLAB®):	Epochs:	Final error:
Algebraic	1	2×10^5	8	1	0
Levenberg-Marquardt	50	5×10^7	150	6	10^{-26}
Resilient Backprop.	150	1×10^7	100	150	0.006