

# Neural Networks

1

## Neural Networks

- Analogy to biological neural systems, the most robust learning systems we know.
- Attempt to understand natural biological systems through computational modeling.
- Massive parallelism allows for computational efficiency.
- Help understand “distributed” nature of neural representations (rather than “localist” representation) that allow robustness and graceful degradation.
- Intelligent behavior as an “emergent” property of large number of simple units rather than from explicitly encoded symbolic rules and algorithms.

2

## Neural Speed Constraints

- Neurons have a “switching time” on the order of a few milliseconds, compared to nanoseconds for current computing hardware.
- However, neural systems can perform complex cognitive tasks (vision, speech understanding) in tenths of a second.
- Only time for performing 100 serial steps in this time frame, compared to orders of magnitude more for current computers.
- Must be exploiting “massive parallelism.”
- Human brain has about  $10^{11}$  neurons with an average of  $10^4$  connections each.

3

## Neural Network Learning

- Learning approach based on modeling adaptation in biological neural systems.
- Perceptron: Initial algorithm for learning simple neural networks (single layer) developed in the 1950’ s.
- Backpropagation: More complex algorithm for learning multi-layer neural networks developed in the 1980’ s.

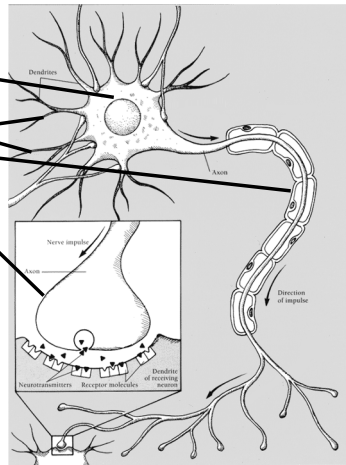
4

4

## Real Neurons

### Cell structures

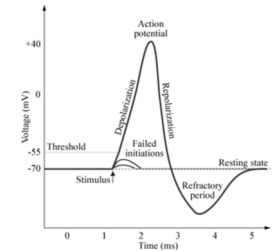
- Cell body
- Dendrites
- Axon
- Synaptic terminals



5

## Neural Communication

- Electrical potential across cell membrane exhibits spikes called action potentials.
- Spike originates in cell body, travels down axon, and causes synaptic terminals to release neurotransmitters.
- Chemical diffuses across synapse to dendrites of other neurons.
- Neurotransmitters can be excitatory or inhibitory.
- If net input of neurotransmitters to a neuron from other neurons is excitatory and exceeds some threshold, it fires an action potential.



6

## Real Neural Learning

- Synapses change size and strength with experience.
- Hebbian learning: When two connected neurons are firing at the same time, the strength of the synapse between them increases.
- “Neurons that fire together, wire together.”

7

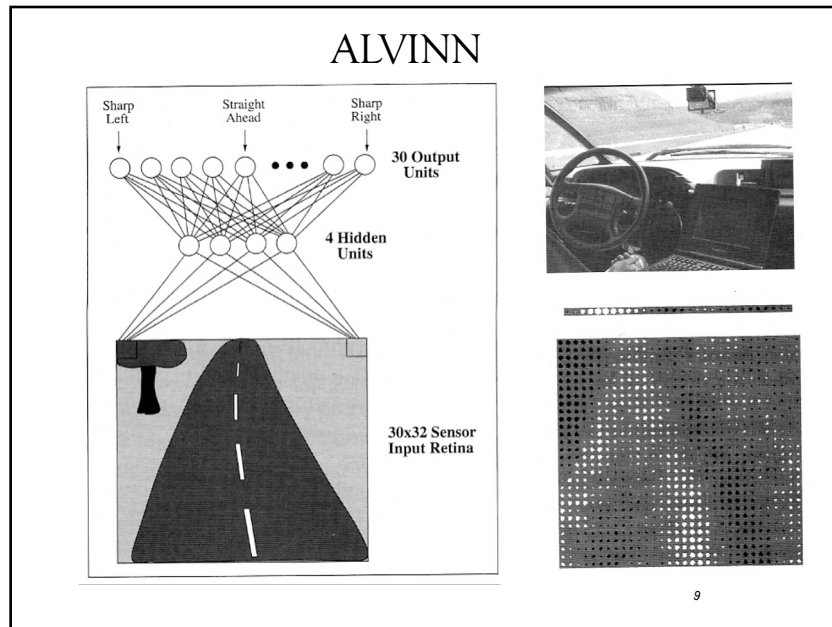
7

## Neural Network Representation

- ALVINN uses a learned ANN to steer an autonomous vehicle driving at normal speeds on public highways
  - Input to network: 30x32 grid of pixel intensities obtained from a forward-pointed camera mounted on the vehicle
  - Output: direction in which the vehicle is steered
  - Trained to mimic observed steering commands of a human driving the vehicle for approximately 5 minutes

8

8



9

## Appropriate problems

- ANN learning well-suited to problems which the training data corresponds to noisy, complex data (inputs from cameras or microphones)
- Can also be used for problems with symbolic representations
- Most appropriate for problems where
  - Instances have many attribute-value pairs
  - Target function output may be discrete-valued, real-valued, or a vector of several real- or discrete-valued attributes
  - Training examples may contain errors
  - Long training times are acceptable
  - Fast evaluation of the learned target function may be required
  - The ability for humans to understand the learned target function is not important

10

10

## Artificial Neurons (1)

- Artificial neurons are based on biological neurons.
- Each neuron in the network receives one or more inputs.
- An activation function is applied to the inputs, which determines the output of the neuron – the activation level.

• *The charts on the right show three typical activation functions.*

11

11

## Artificial Neurons (2)

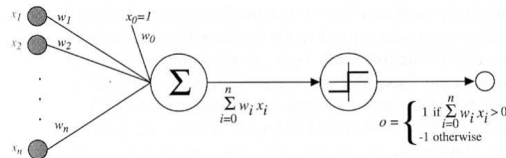
- A typical activation function works as follows:
 
$$X = \sum_{i=1}^n w_i x_i \quad Y = \begin{cases} +1 & \text{for } X > t \\ 0 & \text{for } X \leq t \end{cases}$$
- Each node  $i$  has a weight,  $w_i$  associated with it. The input to node  $i$  is  $x_i$ .
- $t$  is the threshold.
- So if the weighted sum of the inputs to the neuron is above the threshold, then the neuron fires.

12

12

## Perceptrons

- A perceptron is a single neuron that classifies a set of inputs into one of two categories (usually 1 or -1).
- If the inputs are in the form of a grid, a perceptron can be used to recognize visual images of shapes.
- The perceptron usually uses a step function, which returns 1 if the weighted sum of inputs exceeds a threshold, and 0 otherwise.



13

## Artificial Neuron Model

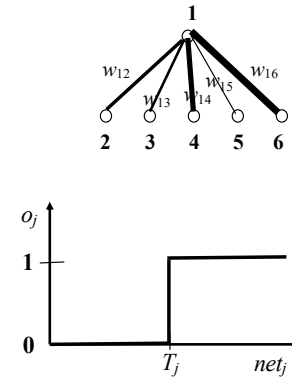
- Model network as a graph with cells as nodes and synaptic connections as weighted edges from node  $i$  to node  $j$ ,  $w_{ji}$
- Model net input to cell as

$$net_j = \sum_i w_{ji} o_i$$

- Cell output is:

$$o_j = \begin{cases} 0 & \text{if } net_j < T_j \\ 1 & \text{if } net_j \geq T_j \end{cases}$$

( $T_j$  is threshold for unit  $j$ )



14

## Neural Computation

- McCollough and Pitts (1943) showed how such model neurons could compute logical functions and be used to construct finite-state machines.
- Can be used to simulate logic gates:
  - AND: Let all  $w_{ji}$  be  $T_j/n$ , where  $n$  is the number of inputs.
  - OR: Let all  $w_{ji}$  be  $T_j$
  - NOT: Let threshold be 0, single input with a negative weight.
- Can build arbitrary logic circuits, sequential machines, and computers with such gates.
- Given negated inputs, two layer network can compute any boolean function using a two level AND-OR network.

15

15

## Perceptron Training

- Assume supervised training examples giving the desired output for a unit given a set of known input activations.
- Learn synaptic weights so that unit produces the correct output for each example.
- Perceptron uses iterative update algorithm to learn a correct set of weights.

16

16

## Gradient Descent Algorithm

Repeat until convergence {

$$\theta_j \leftarrow \theta_j - \eta \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \text{ (for } j=0 \text{ and } j=1)$$

}

Correct simultaneous updates

$$\begin{aligned} \text{temp0} &\leftarrow \theta_0 - \eta \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) \\ \text{temp1} &\leftarrow \theta_1 - \eta \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \\ \theta_0 &\leftarrow \text{temp0} \\ \theta_1 &\leftarrow \text{temp1} \end{aligned}$$

17

17

## Impact of Learning Rate

For simplicity Consider

$$\theta_1 \leftarrow \theta_1 - \eta \frac{\partial}{\partial \theta_1} J(\theta_1)$$

If the learning rate is too small, the gradient descent can be too slow.

If the learning rate is too large, the gradient descent can overshoot the minimum. It may fail to converge, or even diverge.

18

18

## Perceptron Learning Rule

○ Update weights by:

$$w_{ji} = w_{ji} + \eta(t_j - o_j)o_i$$

where  $\eta$  is the “learning rate”

$t_j$  is the teacher specified output for unit  $j$ .

○ Equivalent to rules:

- If output is correct do nothing.
- If output is high, lower weights on active inputs
- If output is low, increase weights on active inputs

○ Also adjust threshold to compensate:

$$T_j = T_j - \eta(t_j - o_j)$$

19

19

## Perceptron Learning Algorithm

○ Iteratively update weights until convergence.

Initialize weights to random values

Until outputs of all training examples are correct

For each training pair,  $E$ , do:

Compute current output  $o_j$  for  $E$  given its inputs

Compare current output to target value,  $t_j$ , for  $E$

Update synaptic weights and threshold using learning rule

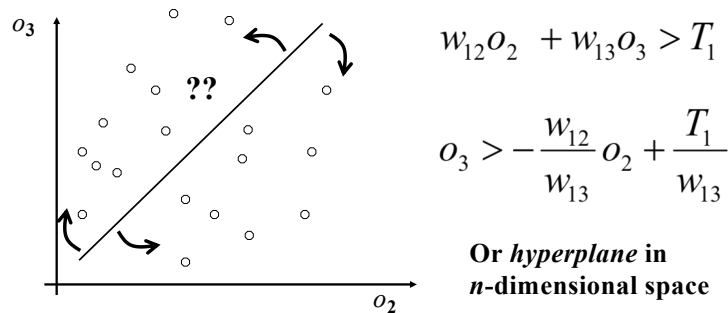
○ Each execution of the outer loop is typically called an *epoch*.

20

20

## Perceptron as a Linear Separator

- Since perceptron uses linear threshold function, it is searching for a linear separator that discriminates the classes.

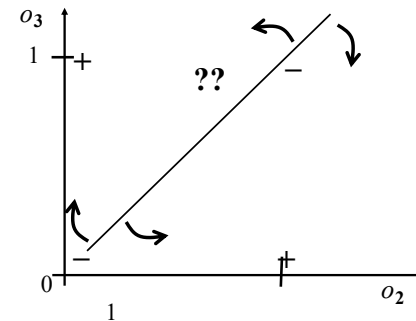


21

21

## Concept Perceptron Cannot Learn

- Cannot learn exclusive-or, or parity function in general.



22

22

## Perceptron Limits

- System obviously cannot learn concepts it cannot represent.
- Minsky and Papert (1969) wrote a book analyzing the perceptron and demonstrating many functions it could not learn.
- These results discouraged further research on neural nets; and symbolic AI became the dominate paradigm.

23

23

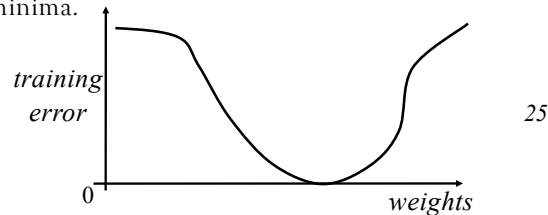
## Perceptron Convergence and Cycling Theorems

- **Perceptron convergence theorem:** If the data is linearly separable and therefore a set of weights exist that are consistent with the data, then the Perceptron algorithm will eventually converge to a consistent set of weights.
- **Perceptron cycling theorem:** If the data is not linearly separable, the Perceptron algorithm will eventually repeat a set of weights and threshold at the end of some epoch and therefore enter an infinite loop.
  - By checking for repeated weights+threshold, one can guarantee termination with either a positive or negative result.

24

## Perceptron as Hill Climbing

- The hypothesis space being search is a set of weights and a threshold.
- Objective is to minimize classification error on the training set.
- Perceptron effectively does hill-climbing (gradient descent) in this space, changing the weights a small amount at each point to decrease training set error.
- For a single model neuron, the space is well behaved with a single minima.



25

## Perceptron Performance

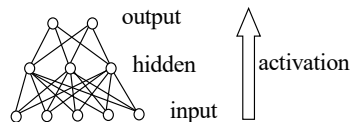
- Linear threshold functions are restrictive (high bias) but still reasonably expressive; more general than:
  - Pure conjunctive
  - Pure disjunctive
  - M-of-N (at least M of a specified set of N features must be present)
- In practice, converges fairly quickly for linearly separable data.
- Can effectively use even incompletely converged results when only a few outliers are misclassified.
- Experimentally, Perceptron does quite well on many benchmark data sets.

26

26

## Multi-Layer Networks

- Multi-layer networks can represent arbitrary functions, but an effective learning algorithm for such networks was thought to be difficult.
- A typical multi-layer network consists of an input, hidden and output layer, each fully connected to the next, with activation feeding forward.



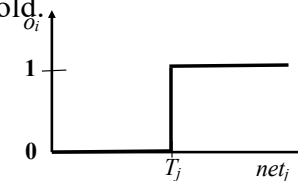
- The weights determine the function computed. Given an arbitrary number of hidden units, any boolean function can be computed with a single hidden layer.

27

27

## Hill-Climbing in Multi-Layer Nets

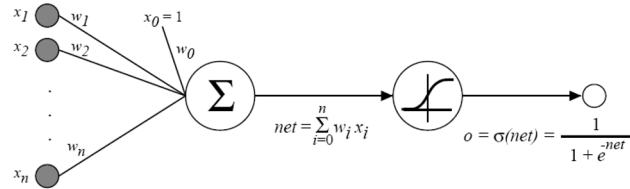
- Since “greed is good” perhaps hill-climbing can be used to learn multi-layer networks in practice although its theoretical limits are clear.
- However, to do gradient descent, we need the output of a unit to be a differentiable function of its input and weights.
- Standard linear threshold function is not differentiable at the threshold.



28

28

## Sigmoid Unit



- $\sigma(x)$  is the sigmoid function  $\frac{1}{1 + e^{-x}}$
- For threshold  $T_j$   $o_j = \frac{1}{1 + e^{-(net_j - T_j)}}$
- Nice property: differentiable  $\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x))$
- Derive gradient descent rules to train
  - One sigmoid unit - node
  - Multilayer networks of sigmoid units

29

39

## Gradient Descent

- Define objective to minimize error:

$$E(W) = \sum_{d \in D} \sum_{k \in K} (t_{kd} - o_{kd})^2$$

where  $D$  is the set of training examples,  $K$  is the set of output units,  $t_{kd}$  and  $o_{kd}$  are, respectively, the teacher and current output for unit  $k$  for example  $d$ .

- The derivative of a sigmoid unit with respect to net input is:

$$\frac{\partial o_j}{\partial net_j} = o_j(1 - o_j)$$

- Learning rule to change weights to minimize error is:

$$\Delta w_{ji} = -\eta \frac{\partial E}{\partial w_{ji}} \quad 30$$

30

## Backpropagation Learning Rule

- Each weight changed by:

$$\Delta w_{ji} = \eta \delta_j o_i$$

$$\delta_j = o_j(1 - o_j)(t_j - o_j) \quad \text{if } j \text{ is an output unit}$$

$$\delta_j = o_j(1 - o_j) \sum_k \delta_k w_{kj} \quad \text{if } j \text{ is a hidden unit}$$

where  $\eta$  is a constant called the learning rate

$t_j$  is the correct teacher output for unit  $j$

$\delta_j$  is the error measure for unit  $j$

31

31

## Back propagation Algorithm

- Create a feed-forward network with  $n_{in}$  inputs,  $n_{hidden}$  hidden units, and  $n_{out}$  output units.
- Initialize all network weights to small random numbers
- Until termination condition is met, Do
  - For each  $\langle x, t \rangle$  in training examples, Do

Propagate the input forward through the network:

1. Input the instance  $x$  to the network and compute the output  $o_u$  of every unit  $u$  in the network

Propagate the errors backward through the network:

2. For each network output unit  $k$ , calculate its error term  $\delta_k$

3. For each hidden unit  $h$ , calculate its error term  $\delta_h$

$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in \text{Outputs}} w_{kh} \delta_k$$

4. Update each network weight  $w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$  where

$$\Delta w_{ji} = \alpha \delta_j x_{ji}$$

32

32



## Error Backpropagation

- First calculate error of output units and use this to change the top layer of weights.

Current output:  $o_j=0.2$

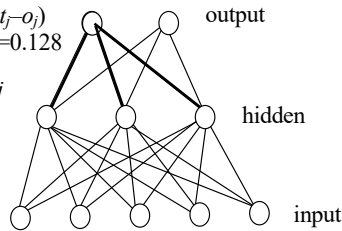
Correct output:  $t_j=1.0$

Error  $\delta_j = o_j(1-o_j)(t_j-o_j)$

$0.2(1-0.2)(1-0.2)=0.128$

Update weights into  $j$

$$\Delta w_{ji} = \eta \delta_j o_i$$



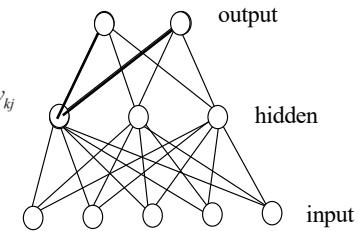
33

33

## Error Backpropagation

- Next calculate error for hidden units based on errors on the output units it feeds into.

$$\delta_j = o_j(1-o_j) \sum_k \delta_k w_{kj}$$



34

34

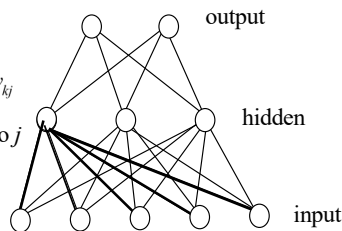
## Error Backpropagation

- Finally update bottom layer of weights based on errors calculated for hidden units.

$$\delta_j = o_j(1-o_j) \sum_k \delta_k w_{kj}$$

Update weights into  $j$

$$\Delta w_{ji} = \eta \delta_j o_i$$



35

35

## Backpropagation Training Algorithm

Create the 3-layer network with  $H$  hidden units with full connectivity between layers.

Set weights to small random real values.

Until all training examples produce the correct value (within  $\epsilon$ ), or mean squared error ceases to decrease, or other termination criteria:

Begin epoch

For each training example,  $d$ , do:

Calculate network output for  $d$ 's input values

Compute error between current output and correct output for  $d$

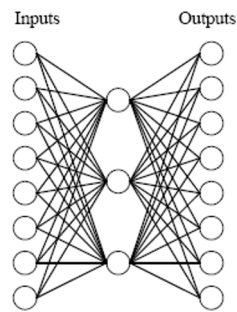
Update weights by back propagating error and using learning rule

End epoch

36

36

## Hidden Layer representation



Target Function:

Input	Output
10000000	→ 10000000
01000000	→ 01000000
00100000	→ 00100000
00010000	→ 00010000
00001000	→ 00001000
00000100	→ 00000100
00000010	→ 00000010
00000001	→ 00000001

*Can this be learned?*

37

37

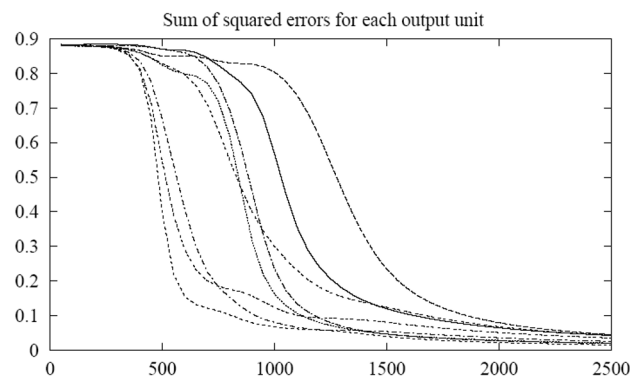
Yes

Input	Hidden Values	Output
10000000	→ .89 .04 .08	→ 10000000
01000000	→ .15 .99 .99	→ 01000000
00100000	→ .01 .97 .27	→ 00100000
00010000	→ .99 .97 .71	→ 00010000
00001000	→ .03 .05 .02	→ 00001000
00000100	→ .01 .11 .88	→ 00000100
00000010	→ .80 .01 .98	→ 00000010
00000001	→ .60 .94 .01	→ 00000001

38

38

## Plots of Squared Error



39

39

## Comments on Training Algorithm

- Not guaranteed to converge to zero training error, may converge to local optima or oscillate indefinitely.
- However, in practice, does converge to low error for many large networks on real data.
- Many epochs (thousands) may be required, hours or days of training for large networks.
- To avoid local-minima problems, run several trials starting with different random weights (*random restarts*).
  - Take results of trial with lowest training set error.
  - Build a committee of results from multiple trials (possibly weighting votes by training set accuracy).

40

40

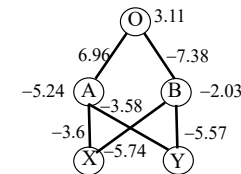
## Representational Power

- **Boolean functions:** Any boolean function can be represented by a two-layer network with sufficient hidden units.
- **Continuous functions:** Any bounded continuous function can be approximated with arbitrarily small error by a two-layer network.
  - Sigmoid functions can act as a set of basis functions for composing more complex functions, like sine waves in Fourier analysis.
- **Arbitrary function:** Any function can be approximated to arbitrary accuracy by a three-layer network.

41

41

## Sample Learned XOR Network



Hidden Unit A represents:  $\neg(X \wedge Y)$   
 Hidden Unit B represents:  $\neg(X \vee Y)$   
 Output O represents:  $A \wedge \neg B = \neg(X \wedge Y) \wedge (X \vee Y)$   
 $= X \oplus Y$

42

42

## Hidden Unit Representations

- Trained hidden units can be seen as newly constructed features that make the target concept linearly separable in the transformed space.
- On many real domains, hidden units can be interpreted as representing meaningful features such as vowel detectors or edge detectors, etc..
- However, the hidden layer can also become a distributed representation of the input  $\mathbf{in}_3$  which each individual unit is not easily interpretable as a meaningful feature.

43

## Successful Applications

- Text to Speech (NetTalk)
- Fraud detection
- Financial Applications
  - HNC (eventually bought by Fair Isaac)
- Chemical Plant Control
  - Pavillion Technologies
- Automated Vehicles
- Game Playing
  - Neurogammon
- Handwriting recognition

44

44