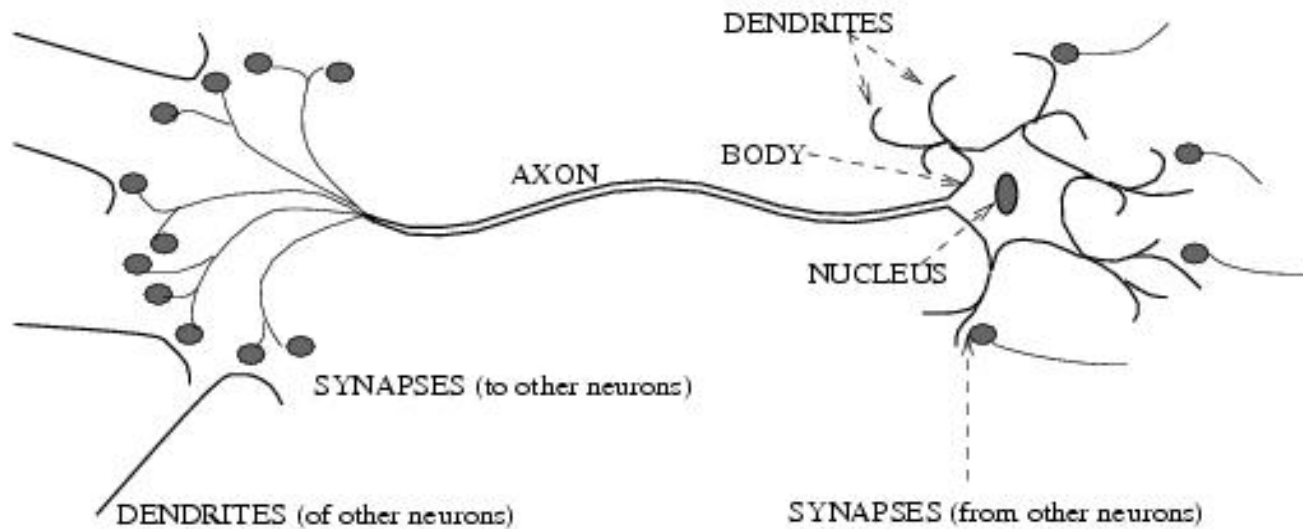


Neural Networks

Biological neural activity



- Each neuron has a *body*, an *axon*, and many *dendrites*
 - Can be in one of the two states: *firing* and *rest*.
 - Neuron fires if total incoming stimulus exceeds a threshold
- *Synapse*: thin gap between axon of one neuron and dendrite of another.
 - Signal exchange
 - Synaptic strength/efficiency

Artificial neural network

- Set of **nodes** (units, neurons, processing elements)
 - Each node has input and output
 - Each node performs a simple computation by its **node function**
- **Weighted connections** between nodes
 - Connectivity gives the structure/architecture of the net
 - What can be computed by a NN is primarily determined by the connections and their weights
- Simplified version of networks of neurons in animal nerve systems

History of NN

- **Pitts & McCulloch (1943)**

- First mathematical model of biological neurons
- All Boolean operations can be implemented by these neuron-like nodes
- Competitor to Von Neumann model for general purpose computing device
- Origin of automata theory

- **Hebb (1949)**

- Hebbian rule of learning: increase the connection strength between neurons i and j whenever both i and j are activated.
- Or increase the connection strength between nodes i and j whenever both nodes are simultaneously ON or OFF.

History: Early booming (50s – early 60s)

– Rosenblatt (1958)

- Perceptron: network of threshold nodes for pattern classification

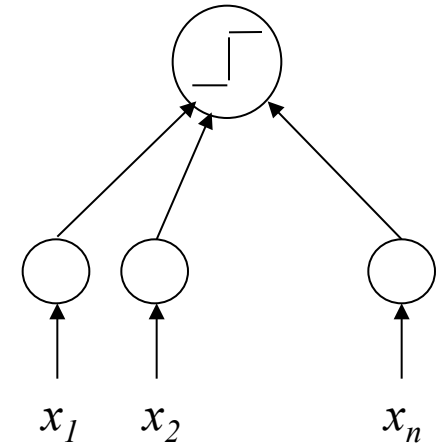
Perceptron learning rule

- Perceptron convergence theorem:
everything that can be represented by a perceptron can be learned

– Widrow and Hoff (1960, 1962)

- Learning rule based on gradient descent (with differentiable unit)

– Minsky's attempt to build a general purpose machine with Pitts/McCulloch units



History: setback in mid 60s – late 70s)

- Serious problems with perceptron model (Minsky's book 1969)
 - Single layer perceptrons cannot represent (learn) simple functions such as XOR
 - Multi-layer of non-linear units may have greater power but there is no learning rule for such nets
 - Scaling problem: connection weights may grow infinitely
 - The first two problems overcame by latter effort in 80's, but the scaling problem persists
- Death of Rosenblatt (1964)
- Striving of Von Neumann machine and AI

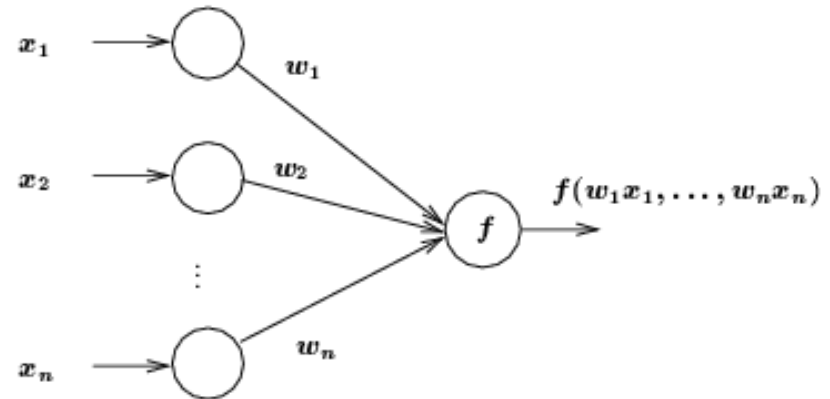
History of NN: Renewed enthusiasm

- New techniques
 - Backpropagation learning for multi-layer feed forward nets (with non-linear, differentiable node functions)
 - Thermodynamic models (Hopfield net, Boltzmann machine, etc.)
 - Unsupervised learning
- Impressive application (character recognition, speech recognition, text-to-speech transformation, process control, associative memory, etc.)
- Traditional approaches face difficult challenges
- Caution:
 - Don't underestimate difficulties and limitations
 - Poses more problems than solutions

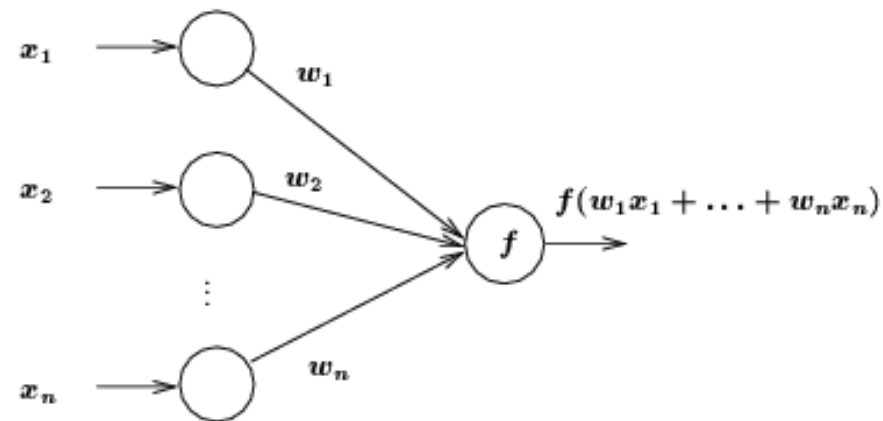
ANN Neuron Models

- Each node has one or more inputs from other nodes, and one output to other nodes
- Input/output values can be
 - Binary $\{0, 1\}$
 - Bipolar $\{-1, 1\}$
 - Continuous (bounded or not)
- All inputs to a node come in at same time and remain activated until output is produced
- Weights associated with links
- Node function

$f(\text{net})$ is the most popular node function where $\text{net} = \sum_{i=1}^n w_i x_i$

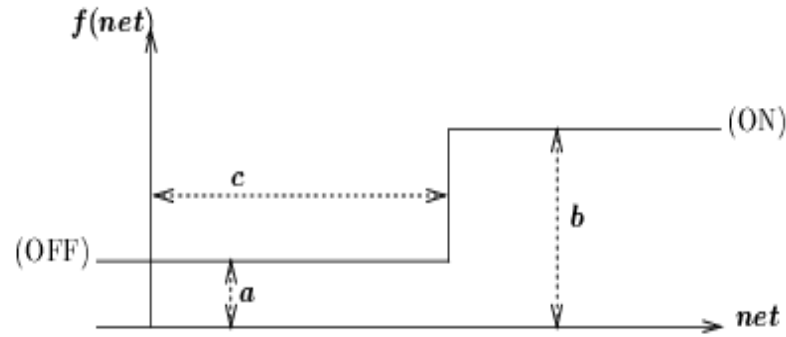


General neuron model

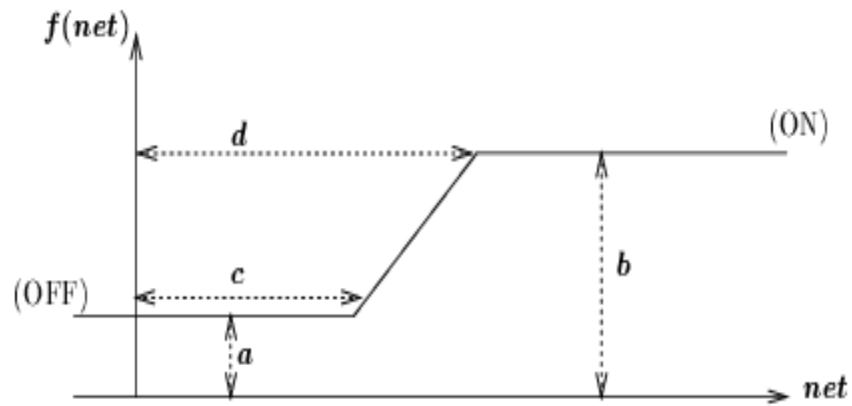


Weighted input summation

Node Function



Step function



Ramp function

Node Function

- **Sigmoid function**

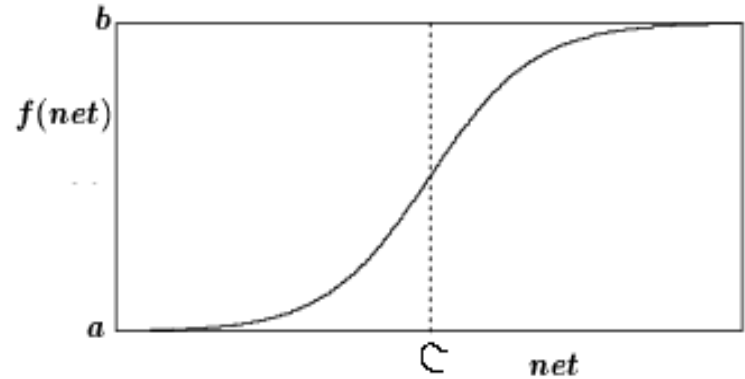
- S-shaped
- Continuous and everywhere differentiable
- Rotationally symmetric about some point ($net = c$)
- Asymptotically approaches saturation points

$$\lim_{net \rightarrow -\infty} f(net) = a \quad \lim_{net \rightarrow \infty} f(net) = b$$

- Examples:

$$f(net) = z + \frac{1}{1 + \exp(-x \cdot net + y)}$$

$$f(net) = \tanh(x \cdot net - y) + z,$$



Sigmoid function

When $y = 0$ and $z = 0$:

$$a = 0, b = 1, c = 0.$$

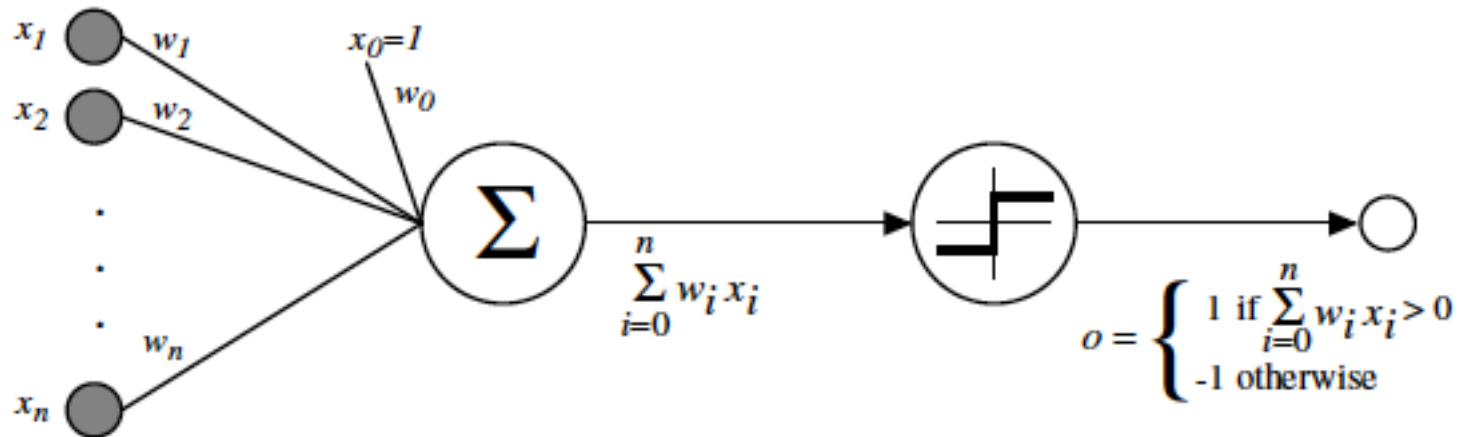
When $y = 0$ and $z = -0.5$

$$a = -0.5, b = 0.5, c = 0.$$

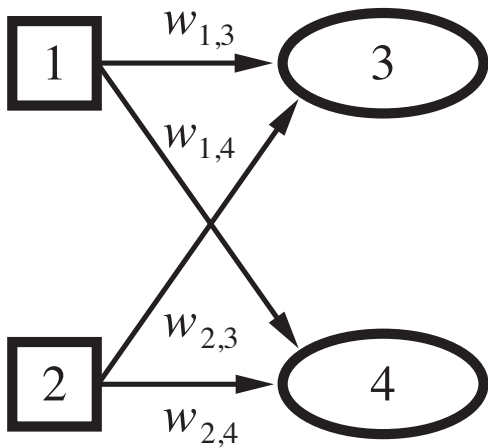
Larger x gives steeper curve

Perceptron

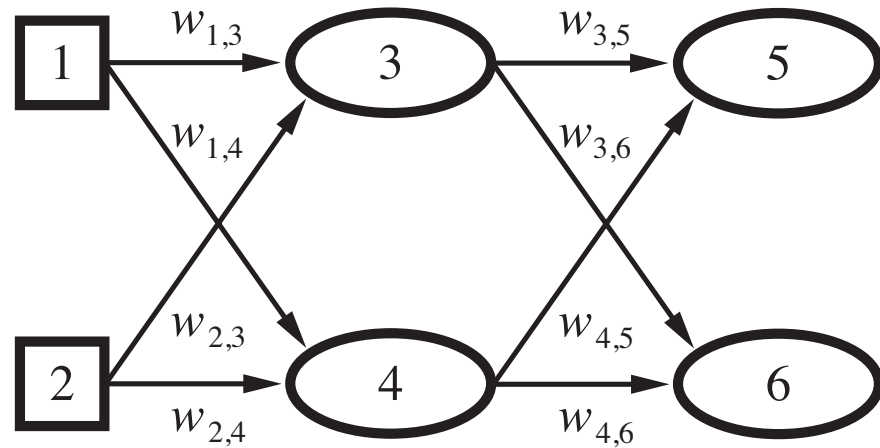
A single layer neural network



Simple architectures



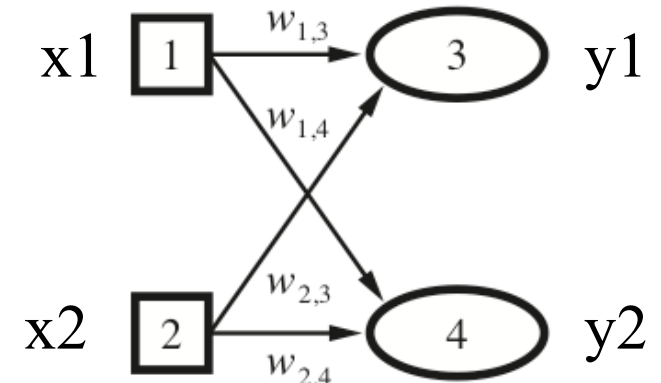
(a)



(b)

Can we make a two bit adder?

- Inputs are bits x_1 and x_2
- Outputs: carry bit (y_1), sum bit (y_2)
- Two NNs, really



| X1 | X2 | Y1 (carry) | Y2 (sum) |
|----|----|------------|----------|
| 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 1 | 1 | 0 |

Perceptron training rule

Adjust weights slightly to reduce error between perceptron output o and target value t ; repeat

$$w_i \leftarrow w_i + \Delta w_i$$

where

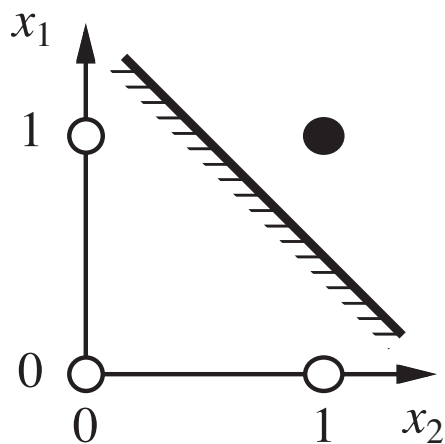
$$\Delta w_i = \eta(t - o)x_i$$

Where:

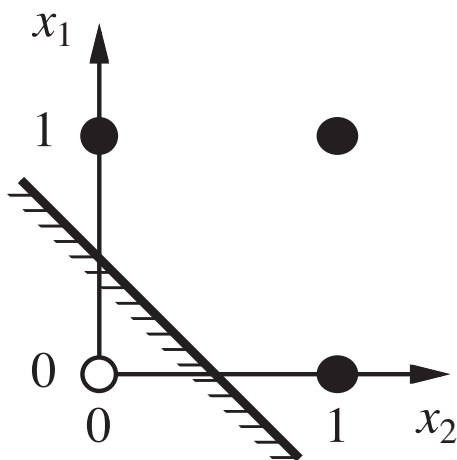
- $t = c(\vec{x})$ is target value
- o is perceptron output
- η is small constant (e.g., .1) called *learning rate*

Not with a perceptron ☹️

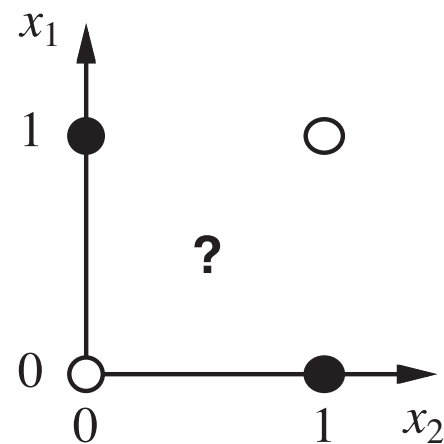
Training examples are not linearly separable for one case: $sum=1$ iff x_1 xor x_2



(a) x_1 and x_2



(b) x_1 or x_2

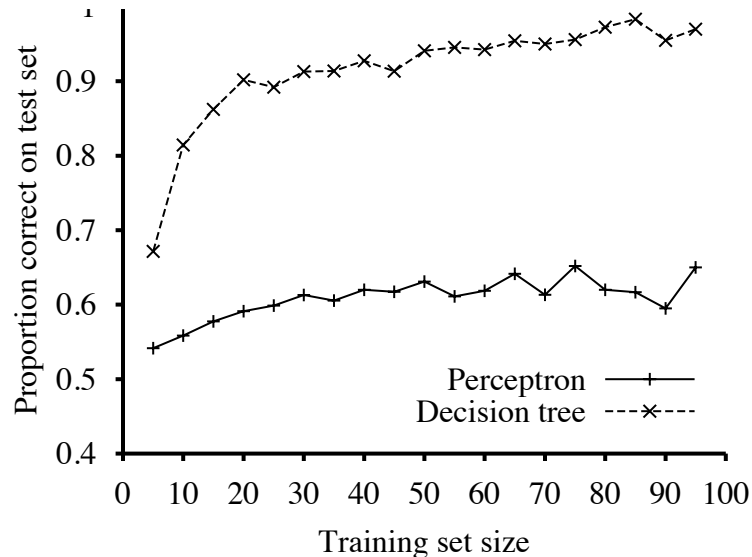
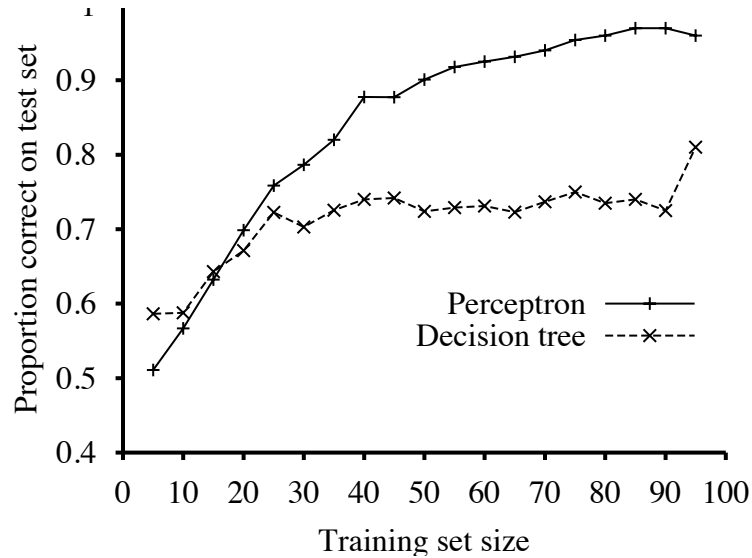


(c) x_1 xor x_2

Works well on some problems

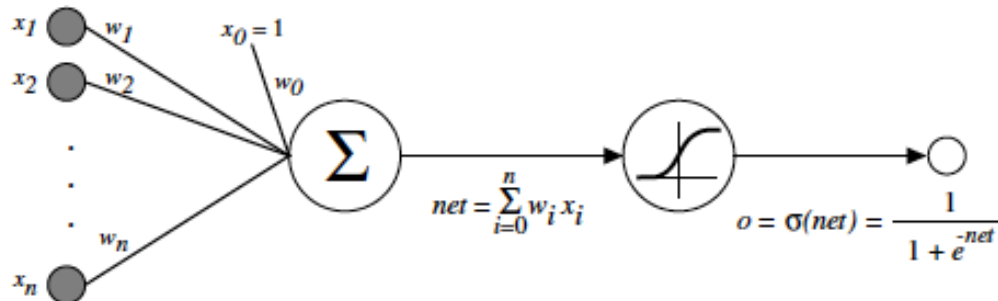
Learning curves

Are majority of inputs 1?



Restaurant example: WillWait?

Sigmoid Unit



$\sigma(x)$ is the sigmoid function

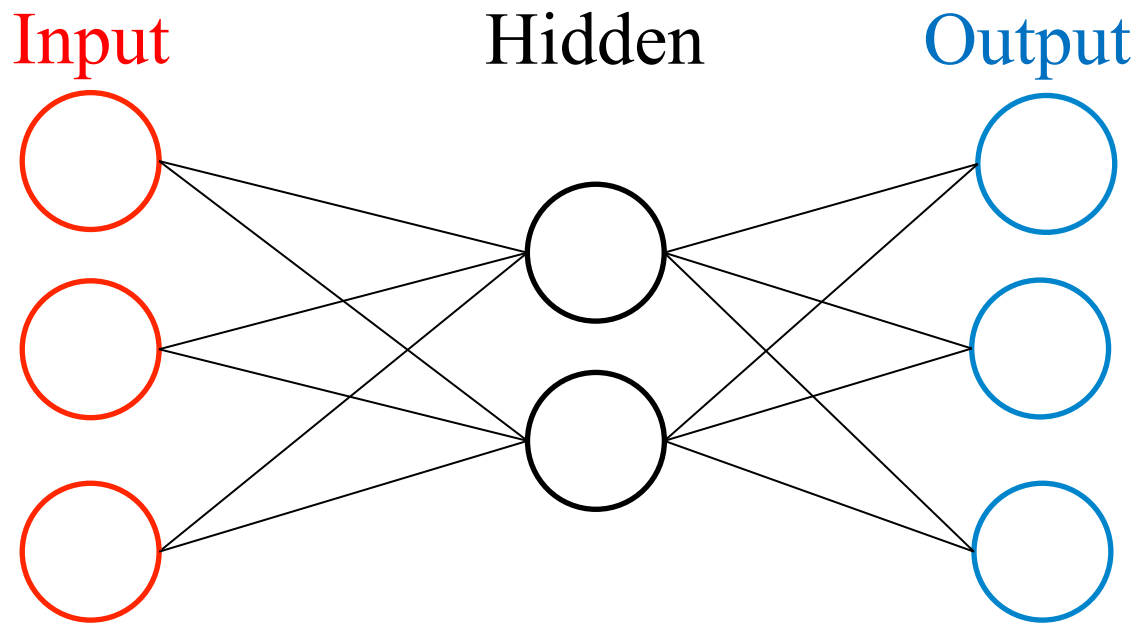
$$\frac{1}{1 + e^{-x}}$$

Nice property: $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$

We can derive gradient decent rules to train

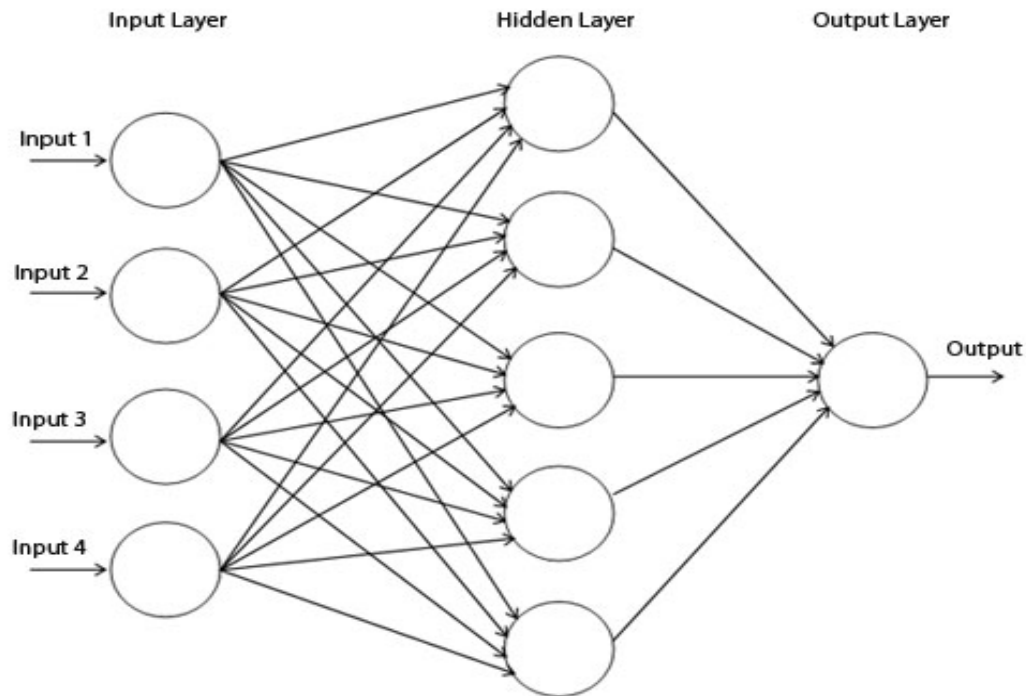
- One sigmoid unit
- *Multilayer networks* of sigmoid units \rightarrow Backpropagation

Multilayer Networks



Backpropagation Algorithm

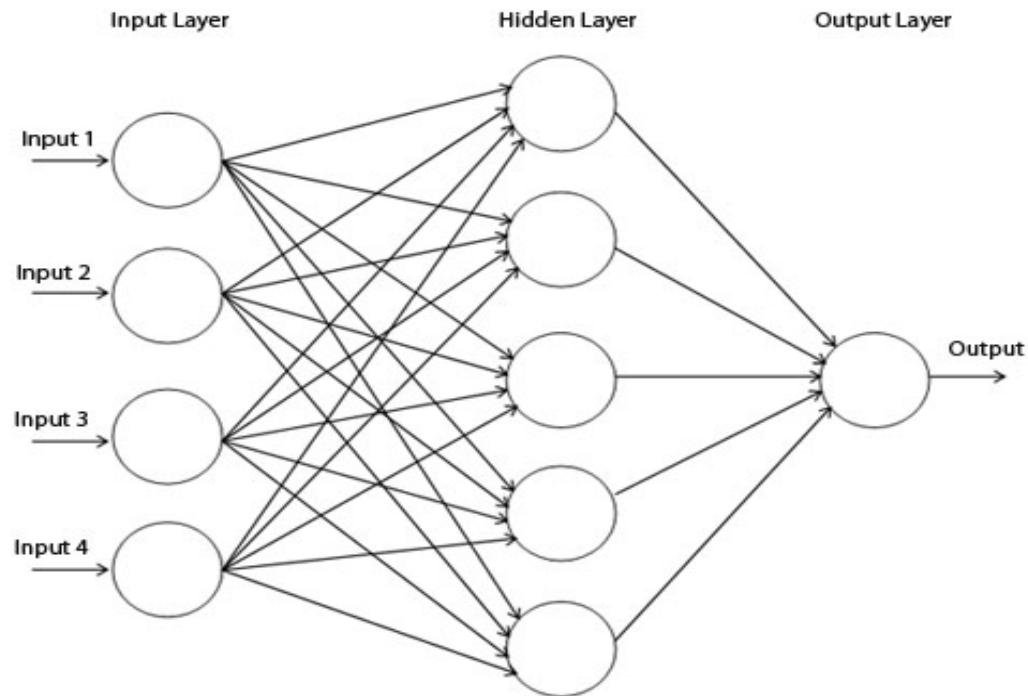
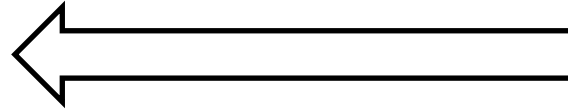
Forward direction
→



Calculate network and error

Backpropagation Algorithm

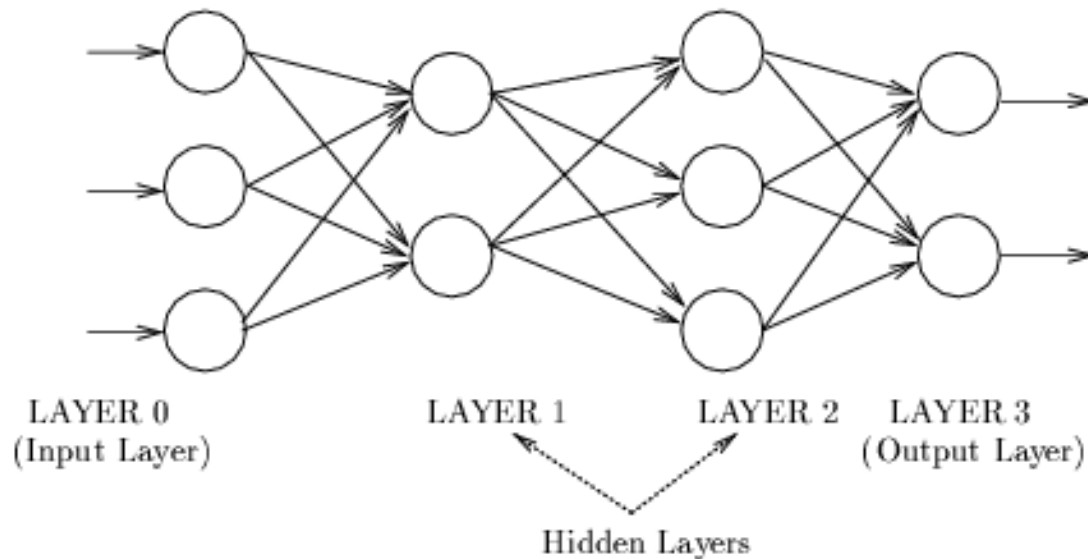
Backward direction



Backpropagate: from output to input, recursively compute $\partial E / \partial w_{ij} = \nabla_{w_{ij}} E$ and adjust weights

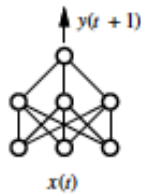
Network Architecture: Feedforward net

- A connection is allowed from a node in layer i only to nodes in layer $i + 1$.
- Most widely used architecture.

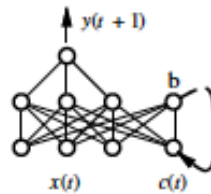


Conceptually, nodes at higher levels successively abstract features from preceding layers

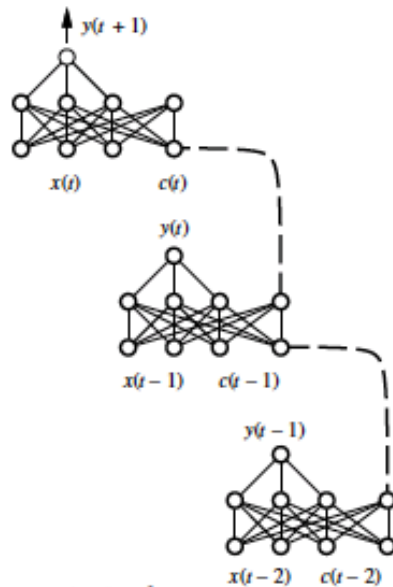
Recurrent neural networks



(a) Feedforward network



(b) Recurrent network



(c) Recurrent network unfolded in time

- Good for learning sequences of data
- e.g., text
- Lots of variations today: convoluted NNs, LSTMs, ...

Neural network playground

The interface displays the following configuration and results:

- Epoch:** 000,000
- Learning rate:** 0.03
- Activation:** Sigmoid
- Regularization:** None
- Regularization rate:** 0
- Problem type:** Classification

DATA: Which dataset do you want to use? (Options: XOR, Spiral, Noise, Sinusoidal)

FEATURES: Which properties do you want to feed in? (Options: X_1 , X_2 , X_1^2 , X_2^2 , $X_1 X_2$, $\sin(X_1)$, $\sin(X_2)$)

2 HIDDEN LAYERS: 4 neurons, 2 neurons

OUTPUT: Test loss 0.554, Training loss 0.543

Legend: Colors show data, neuron and weight values. A color scale from -1 (blue) to 1 (orange) is provided.

Controls: REGENERATE, Show test data, Discretize output

Annotations:
- "The outputs are mixed with varying weights, shown by the thickness of the lines."
- "This is the output from one neuron. Hover to see it larger."