# Lecture 13: Neural Networks
# CS486/686 Intro to Artificial Intelligence

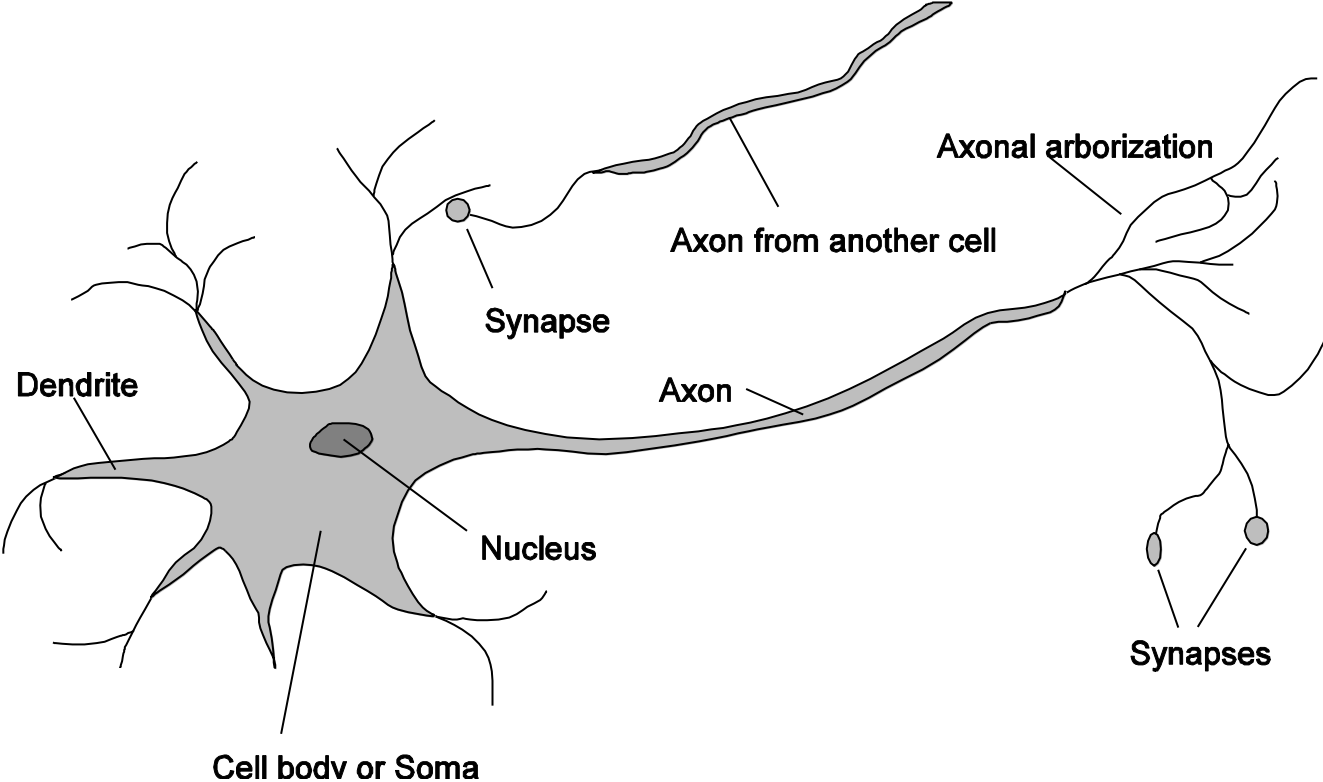Pascal Poupart
David R. Cheriton School of Computer Science

## UNIVERSITY OF
## WATERLOO

# Outline

- Neural networks

  - Perceptron

  - Supervised learning algorithms for neural networks

# Neuron

UNIVERSITY OF
WATERLOO

# Artificial Neural Networks

- Idea: **mimic the brain to do computation**

- Artificial neural network:
  - Nodes (a.k.a. units) correspond to neurons
  - Links correspond to synapses

- Computation:
  - Numerical signal transmitted between nodes corresponds to chemical signals between neurons
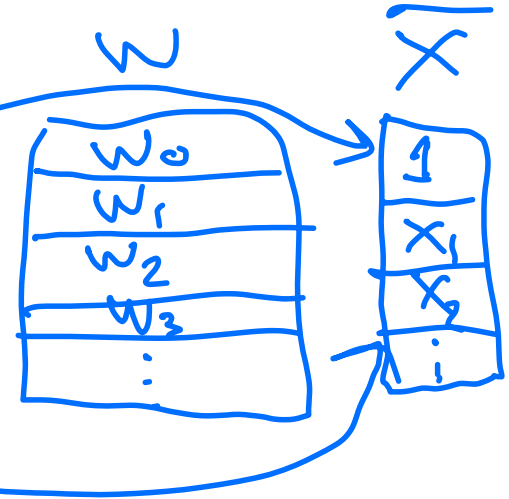  - Nodes modifying numerical signal corresponds to neurons firing rate

UNIVERSITY OF
WATERLOO

# ANN Unit

For each unit i:

- **Weights:** $W$
  - Strength of the link from unit $i$ to unit $j$
  - Input signals $x_i$ weighted by $W_{ji}$ and linearly combined:

$$a_j = \sum_i W_{ji} x_i + W_{j0} = \boldsymbol{W_j} \, \bar{\boldsymbol{x}}$$
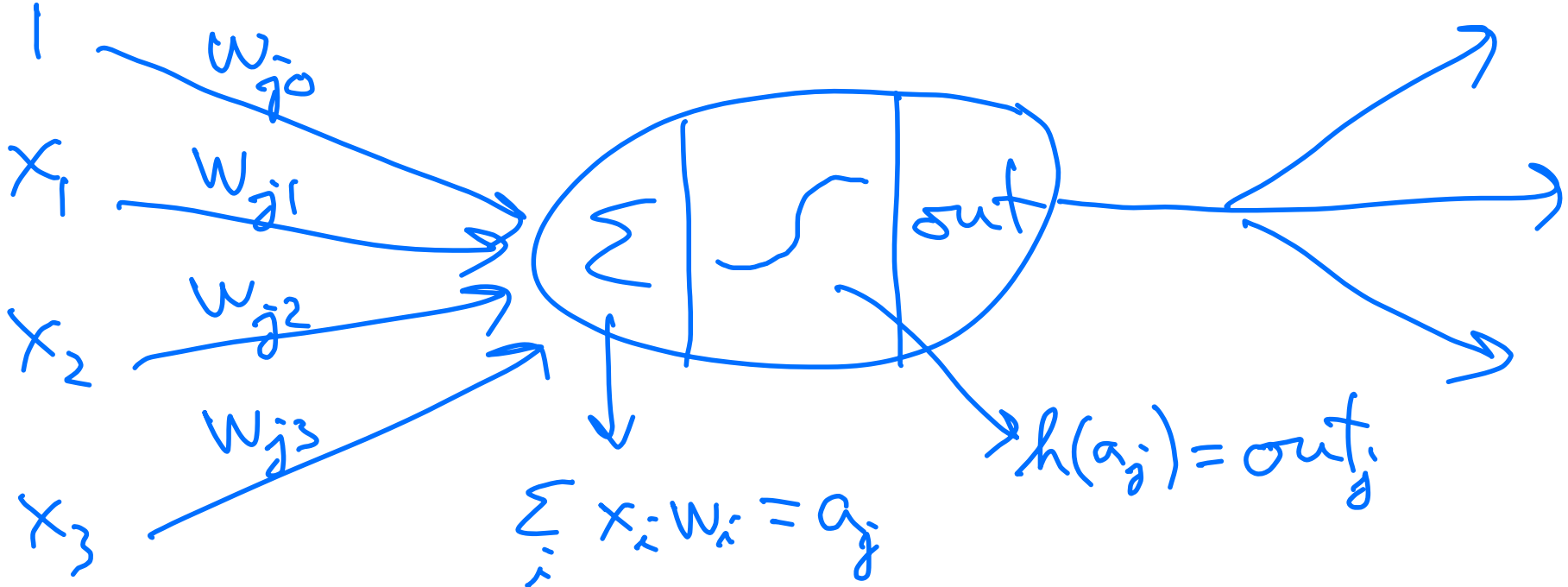


- **Activation function:** $h$
  - Numerical signal produced: $y_j = h(a_j)$

UNIVERSITY OF
WATERLOO

# ANN Unit

- Picture



$$\sum_i x_i \cdot w_i = a_j$$

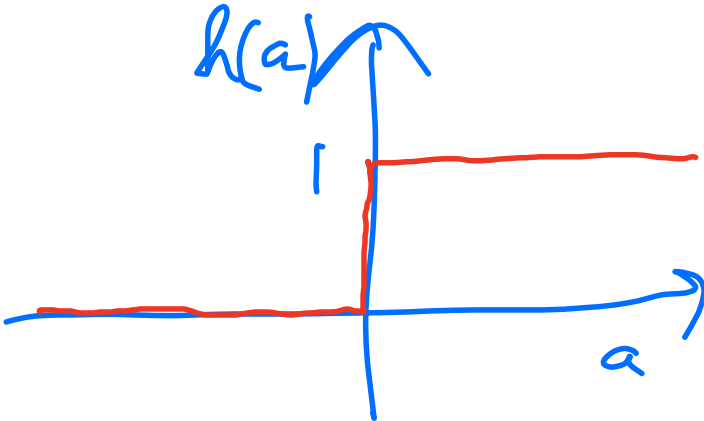$$h(a_j) = out_j$$
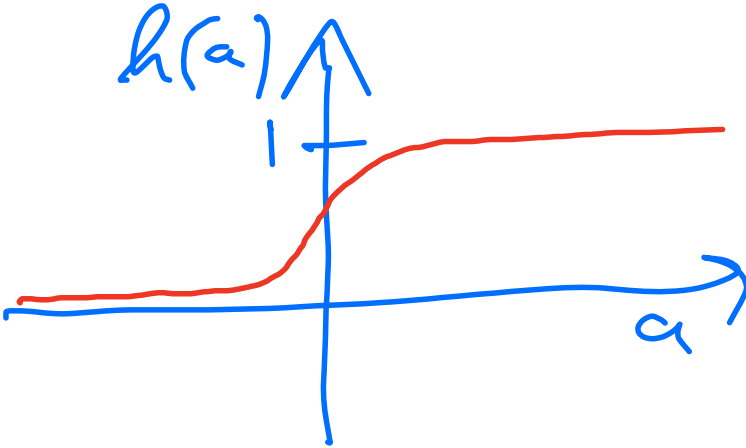
# Activation Function

- Should be nonlinear

  - Otherwise, network is just a linear function

- Often chosen to mimic firing in neurons

  - Unit should be "active" (output near 1) when fed with the "right" inputs

  - Unit should be "inactive" (output near 0) when fed with the "wrong" inputs

UNIVERSITY OF
WATERLOO

# Common Activation Functions

Threshold

Sigmoid

UNIVERSITY OF
WATERLOO

# Logic Gates

$$a = W_0 \cdot 1 + W_1 X_1 + W_2 X_2$$

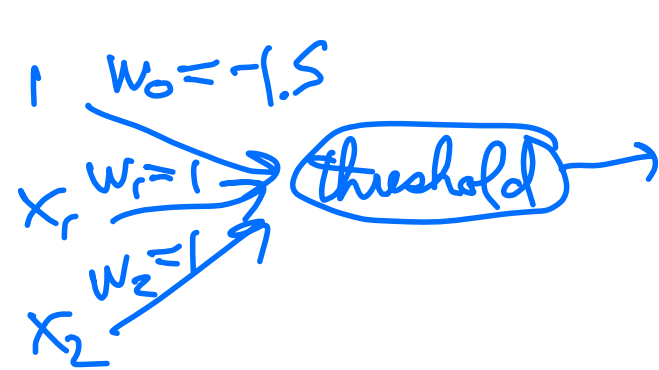$$h(a) = \begin{cases} 0 & \text{if } a \leq 0 \\ 1 & \text{otherwise} \end{cases}$$
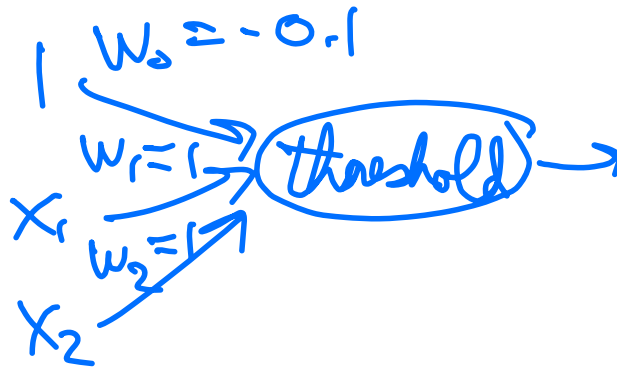
- McCulloch and Pitts (1943)

  - Design ANNs to represent Boolean functions

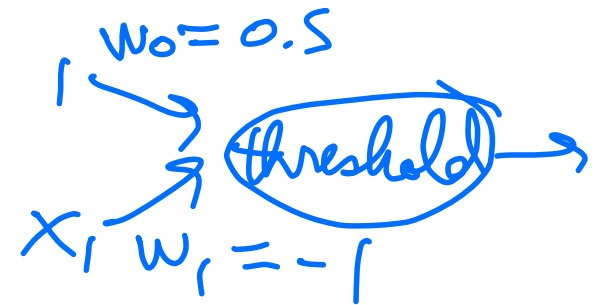- What should be the weights of the following units to code AND, OR, NOT ?

$W_0 = -1.5$
$W_1 = 1$
$X_1$
$W_2 = 1$
$X_2$
threshold

$W_0 = -0.1$
$W_1 = 1$
$X_1$
$W_2 = 1$
$X_2$
threshold

$W_0 = 0.5$
$1$
threshold
$X_1 \quad W_1 = -1$

**AND**                     **OR**                     **NOT**
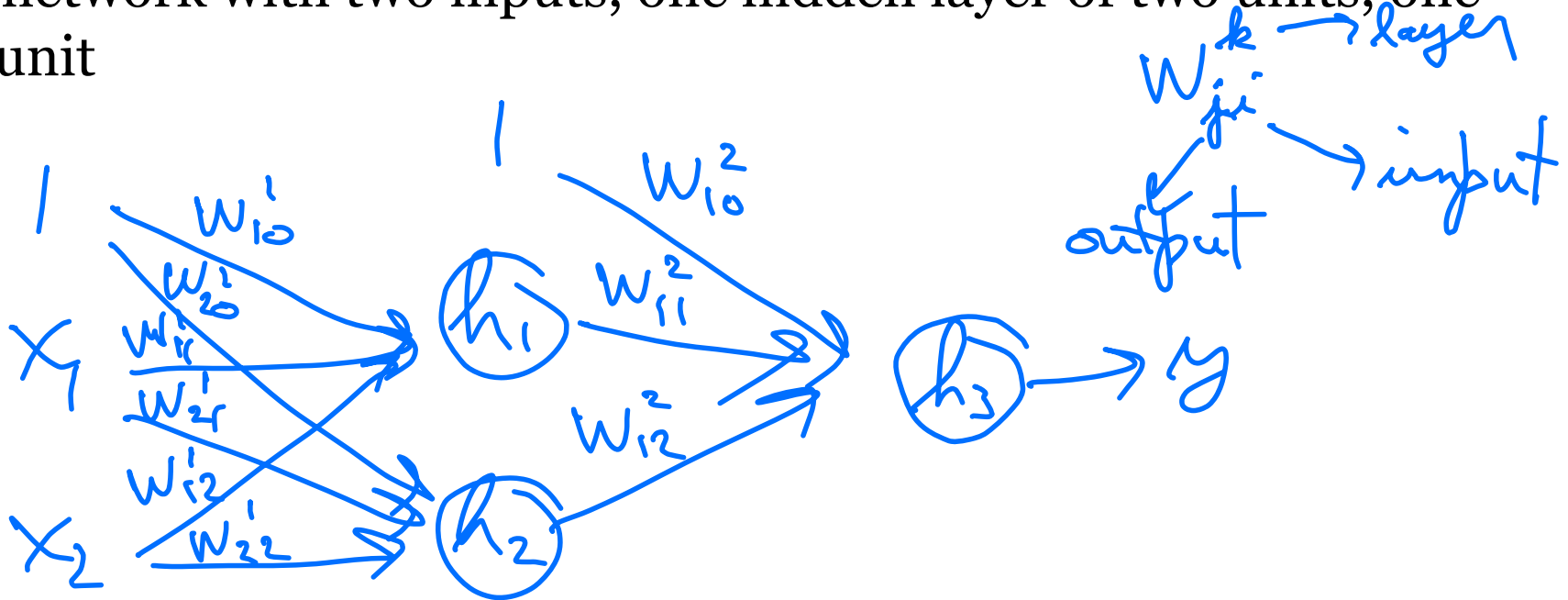
UNIVERSITY OF
WATERLOO

# Network Structures

- **Feed-forward network**
  - Directed **acyclic** graph
  - No internal state
  - Simply computes outputs from inputs


- **Recurrent network**
  - Directed **cyclic** graph
  - Dynamical system with internal states
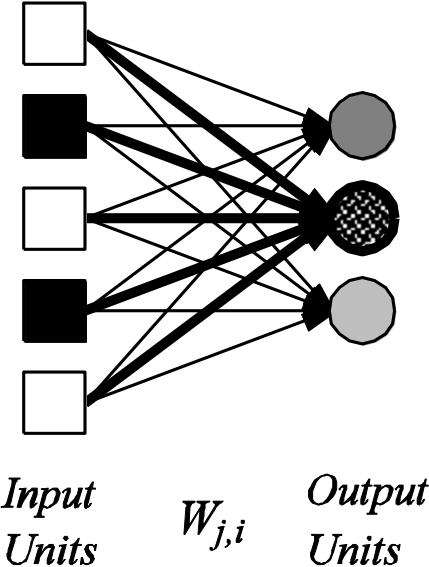  - Can memorize information

UNIVERSITY OF
**WATERLOO**

# Feed-forward network

- Simple network with two inputs, one hidden layer of two units, one output unit

# Perceptron

- Single layer feed-forward network

$$Input \quad W_{j,i} \quad Output$$
$$Units \qquad\qquad Units$$

UNIVERSITY OF
**WATERLOO**

# Threshold Perceptron Hypothesis Space

- Hypothesis space $h_{\boldsymbol{w}}$:
  - All binary classifications with parameters $\boldsymbol{w}$ s.t.
    $$\boldsymbol{w}^T\overline{\boldsymbol{x}} > 0 \rightarrow +1$$
    $$\boldsymbol{w}^T\overline{\boldsymbol{x}} < 0 \rightarrow -1$$

- Since $\boldsymbol{w}^T\overline{\boldsymbol{x}}$ is linear in $\boldsymbol{w}$, perceptron is called a **linear separator**

UNIVERSITY OF
WATERLOO

# Linear Separability

- Are all Boolean gates linearly separable?



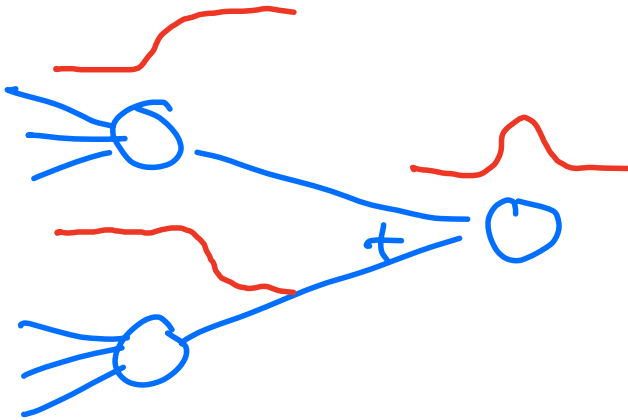(a)   $I_1$ **and** $I_2$     (b)   $I_1$ **or** $I_2$     (c)   $I_1$ **xor** $I_2$

# Sigmoid Perceptron

- Represent "soft" linear separators

Perceptron output

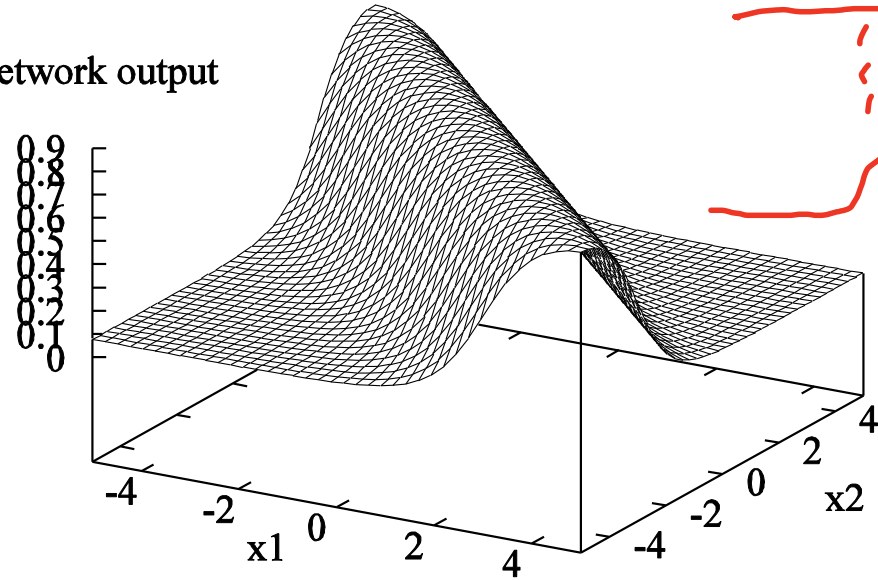# Multilayer Networks

- Adding two sigmoid units with parallel but opposite "cliffs" produces a ridge



Network output

UNIVERSITY OF
WATERLOO

# Multilayer Networks

- Adding two intersecting ridges (and thresholding) produces a bump



Network output

UNIVERSITY OF
WATERLOO
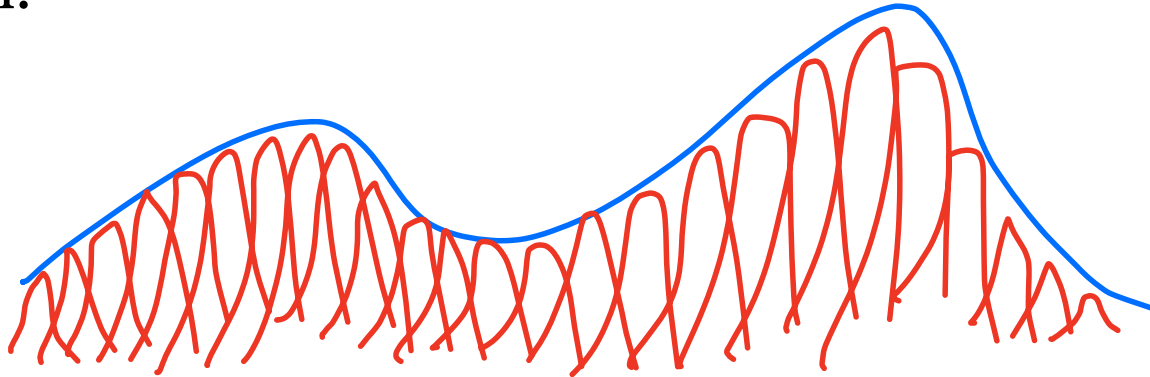
# Multilayer Networks

- By tiling bumps of various heights together, we can approximate any function.



- **Theorem:** Neural networks with at least one hidden layer of sufficiently many sigmoid units can approximate any function arbitrarily closely.

UNIVERSITY OF
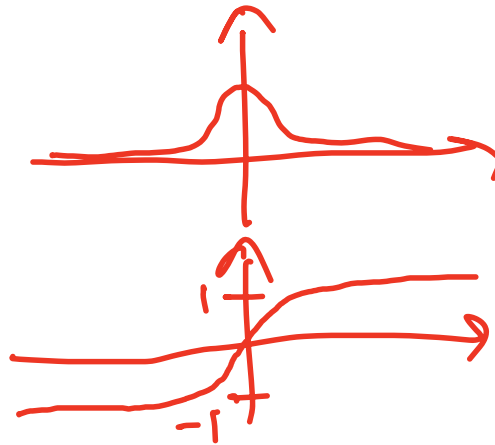**WATERLOO**

# Common activation functions $h$

- Threshold: $h(a) = \begin{cases} 1 & a \geq 0 \\ -1 & a < 0 \end{cases}$

- Sigmoid: $h(a) = \sigma(a) = \frac{1}{1+e^{-a}}$

- Gaussian: $h(a) = e^{-\frac{1}{2}\left(\frac{a-\mu}{\sigma}\right)^2}$

- Tanh: $h(a) = \tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$

- Identity: $h(a) = a$

# Weight training

- Parameters: $< W^{(1)}, W^{(2)}, ... >$

- Objectives:

    - **Error minimization**

        - **Backpropagation (aka "backprop")**

    - Maximum likelihood

    - Maximum a posteriori

    - Bayesian learning

UNIVERSITY OF
**WATERLOO**

# Least squared error

▪ Error function

$$E(\boldsymbol{W}) = \frac{1}{2} \sum_n E_n(\boldsymbol{W})^2 = \frac{1}{2} \sum_n \left\| f(\boldsymbol{x_n}, \boldsymbol{W}) - y_n \right\|_2^2$$

*norm of a vector* → *square*

→ *Euclidean norm*

where $\boldsymbol{x_n}$ is the input of the $n^{th}$ example

$\|a - b\|_2^2 = \sum_i (a_i - b_i)^2$

$y_n$ is the label of the $n^{th}$ example

$f(\boldsymbol{x_n}, \boldsymbol{W})$ is the output of the neural net

UNIVERSITY OF
WATERLOO

# Sequential Gradient Descent

- For each example $(x_n, y_n)$ adjust the weights as follows:

$$w_{ji} \leftarrow w_{ji} - \eta \frac{\partial E_n}{\partial w_{ji}}$$

*(handwritten annotations: $E_m$ curve with $W$ axis; "step size" pointing to $\eta$)*

- How can we compute the gradient efficiently given an arbitrary network structure?

- Answer: **backpropagation algorithm**

- Today: **automatic differentiation**

UNIVERSITY OF
**WATERLOO**

# Backpropagation Algorithm

- Two phases:

  - Forward phase: compute output $z_j$ of each unit $j$



  - Backward phase: compute delta $\delta_j$ at each unit $j$

UNIVERSITY OF
WATERLOO

# Forward phase

- Propagate inputs forward to compute the output of each unit

- Output $z_j$ at unit $j$:

$$z_j = h(a_j) \quad \text{where} \quad a_j = \sum_i w_{ji} z_i$$

UNIVERSITY OF
WATERLOO

# Backward phase

- Use chain rule to recursively compute gradient

    - For each weight $w_{ji}$: $\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}} = \delta_j z_i$

    - Let $\delta_j \equiv \frac{\partial E_n}{\partial a_j}$ then

$$\delta_j = \begin{cases} h'(a_j)(z_j - y_j) & \text{base case: } j \text{ is an output unit} \\ h'(a_j) \sum_k w_{kj} \delta_k & \text{recursion: } j \text{ is a hidden unit} \end{cases}$$

    - Since $a_j = \sum_i w_{ji} z_i$ then $\frac{\partial a_j}{\partial w_{ji}} = z_i$
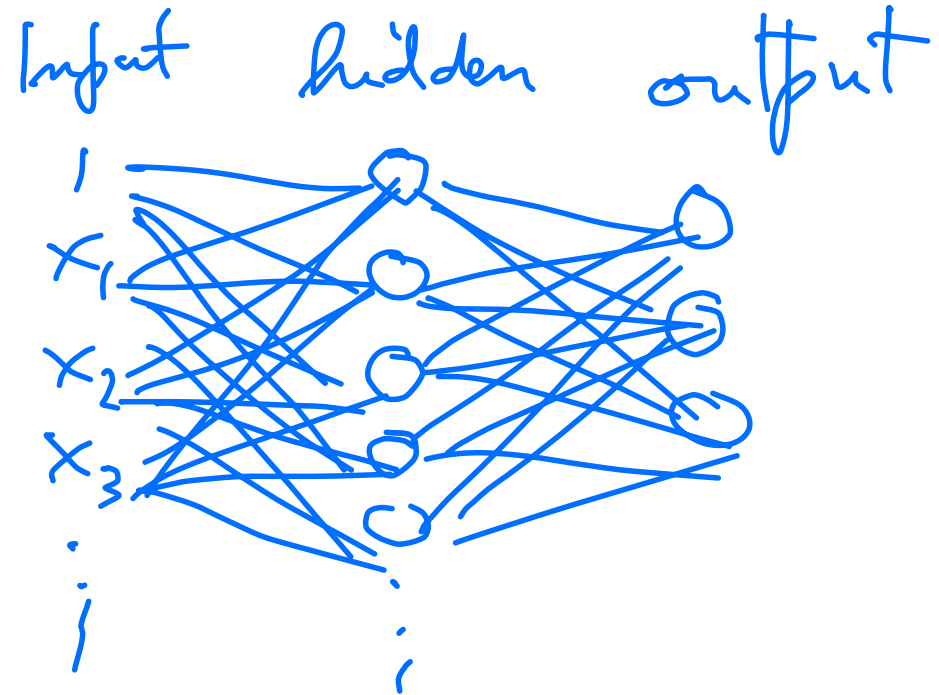
UNIVERSITY OF
WATERLOO

# Simple Example

- Consider a network with two layers:
  - Hidden nodes: $h(a) = \tanh(a) = \dfrac{e^a - e^{-a}}{e^a + e^{-a}}$
    - Tip: $tanh'(a) = 1 - (tanh(a))^2$
  - Output node: $h(a) = a$

- Objective: squared error

UNIVERSITY OF
WATERLOO

# Simple Example

- Forward propagation:
  - Hidden units: $a_j = \sum_i W_{ji} x_i$  $z_j = \tanh(a_j)$
  - Output units: $a_k = \sum_j W_{kj} z_j$  $z_k = a_k$
- Backward propagation:
  - Output units: $\delta_k = z_k - y_k$
  - Hidden units: $\delta_j = (1 - z_j^2) \sum_k W_{kj} \delta_k$
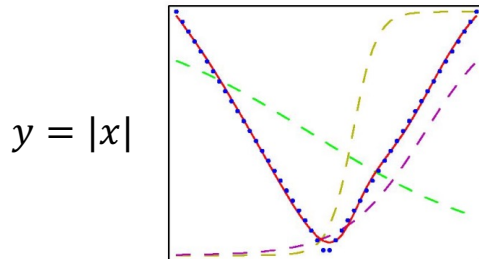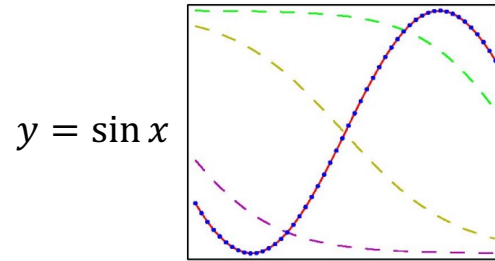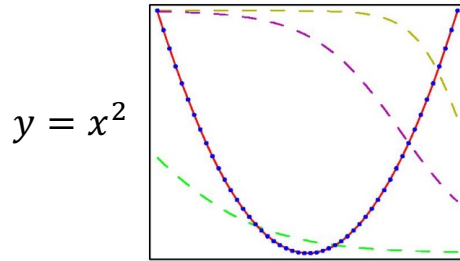- Gradients:
  - Hidden layers: $\frac{\partial E_n}{\partial w_{ji}} = \delta_j x_i = (1 - z_j^2) \sum_k W_{kj} \delta_k x_i$
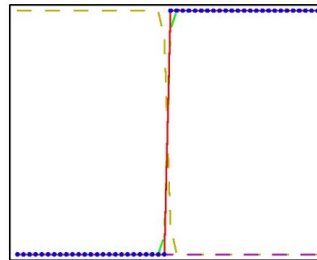  - Output layer: $\frac{\partial E_n}{\partial w_{kj}} = \delta_k z_j = (z_k - y_k) z_j$

UNIVERSITY OF
WATERLOO

# Non-linear regression examples

- Two-layer network:

  - 3 tanh hidden units and 1 identity output unit

$$y = x^2$$

$$y = \sin x$$

$$y = |x|$$

$$y = \int_{-\infty}^{x} \delta(t)dt$$

UNIVERSITY OF
WATERLOO

# Analysis

- Efficiency:

  - Fast gradient computation: linear in number of weights

- Convergence:

  - Slow convergence (linear rate)

  - May get trapped in local optima

- Prone to overfitting

  - Solutions: early stopping, regularization (add $||w||_2^2$ penalty term to objective), dropout

UNIVERSITY OF
WATERLOO