

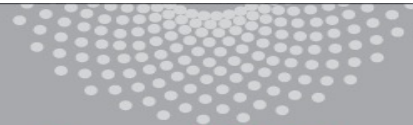
# Neural Networks

Aarti Singh

Machine Learning 10-315  
Feb 14, 2022



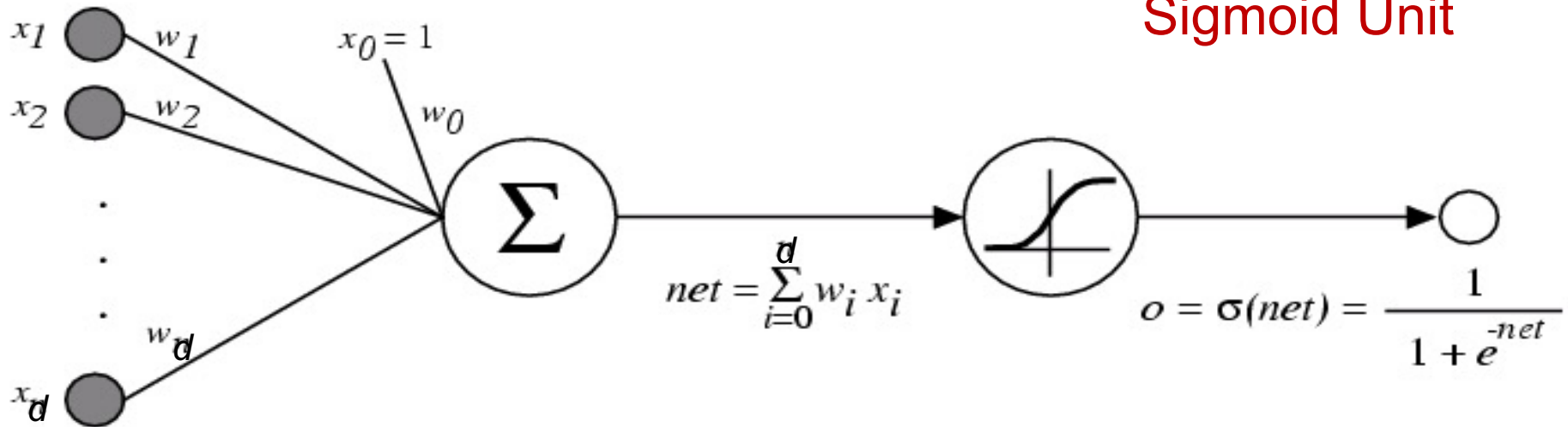
**MACHINE LEARNING** DEPARTMENT



**Carnegie Mellon.**  
School of Computer Science

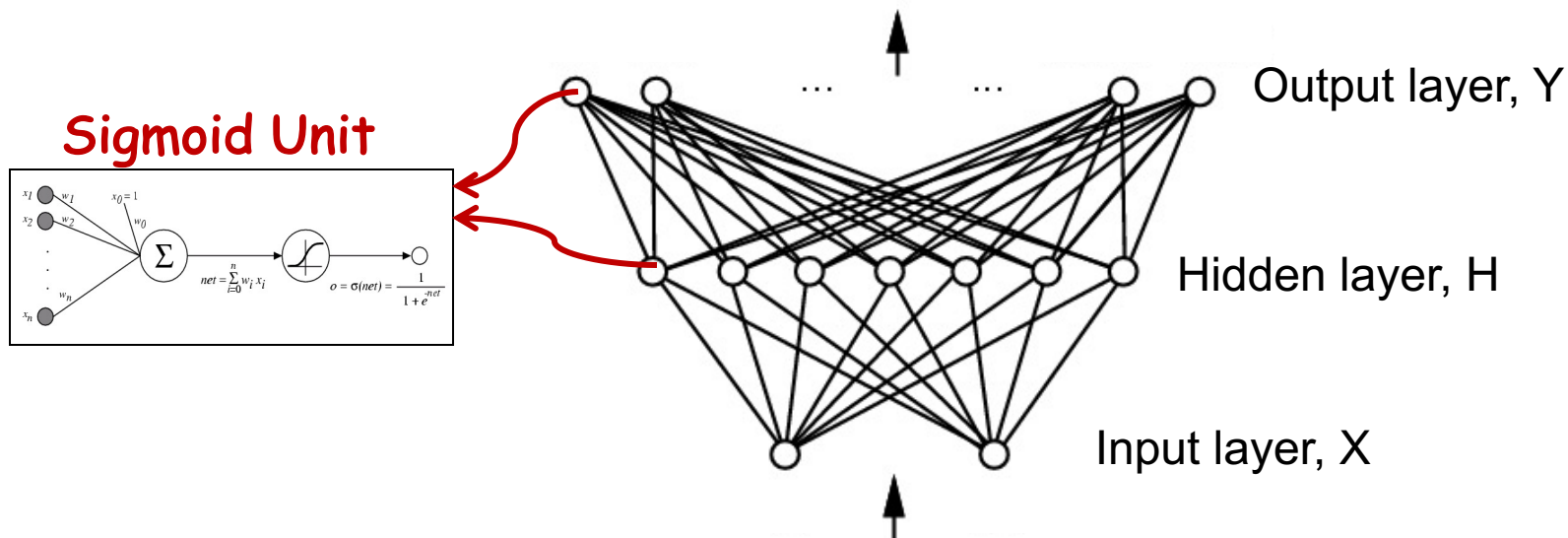
# Logistic function as a Graph

$$\text{Output, } o(\mathbf{x}) = \sigma(w_0 + \sum_i w_i X_i) = \frac{1}{1 + \exp(-(w_0 + \sum_i w_i X_i))}$$



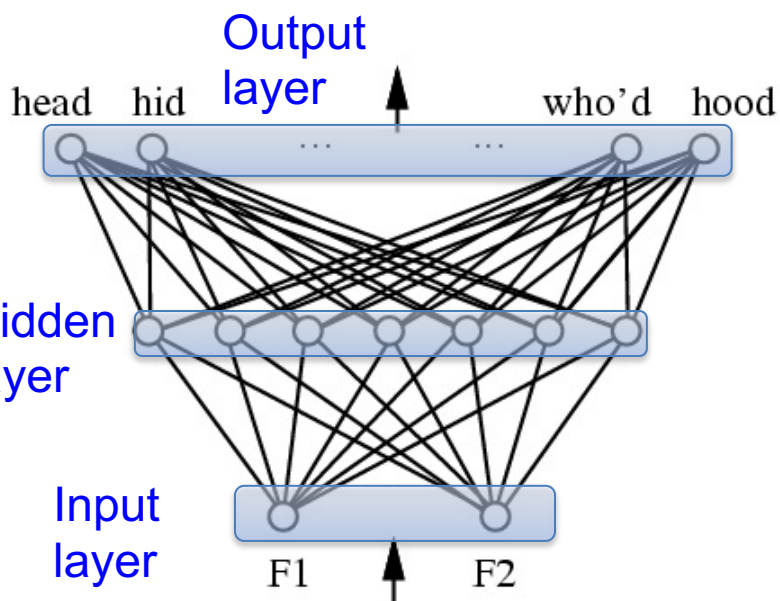
# Neural Networks to learn $f: X \rightarrow Y$

- $f$  can be a **non-linear** function
- $X$  (vector of) continuous and/or discrete variables
- $Y$  (**vector** of) continuous and/or discrete variables
- Neural networks - Represent  $f$  by network of sigmoid (more recently ReLU – next lecture) units :

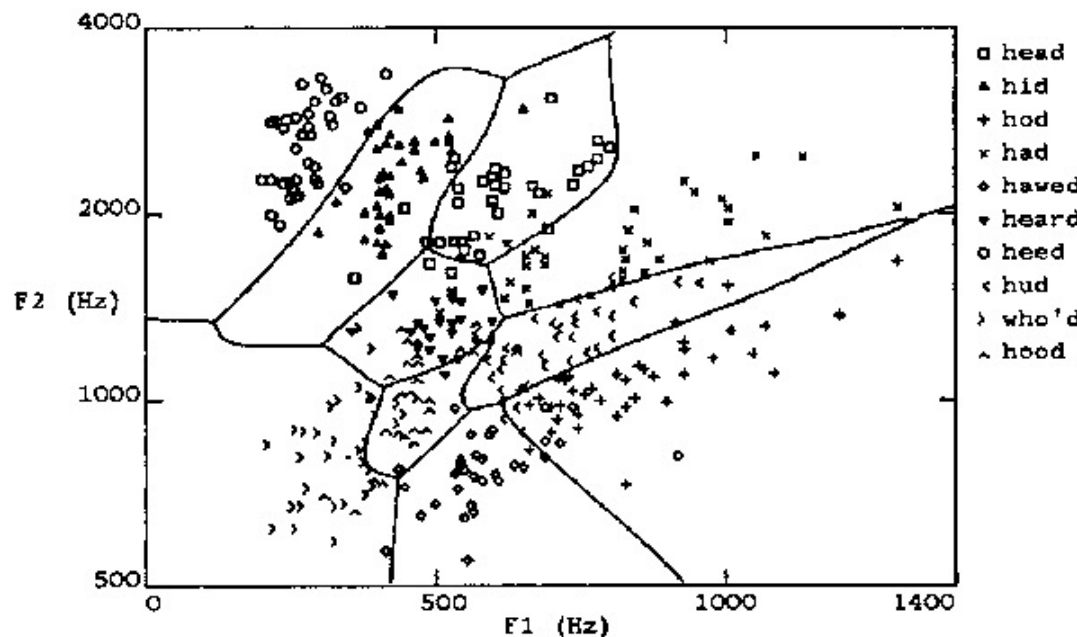


# Multilayer Networks of Sigmoid Units

Neural Network trained to distinguish vowel sounds using 2 formants (features)

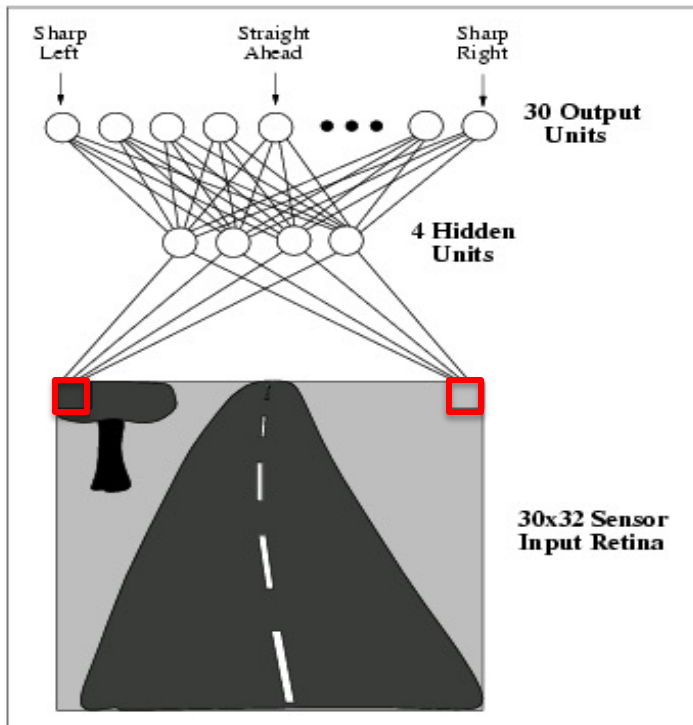


Two layers of logistic units

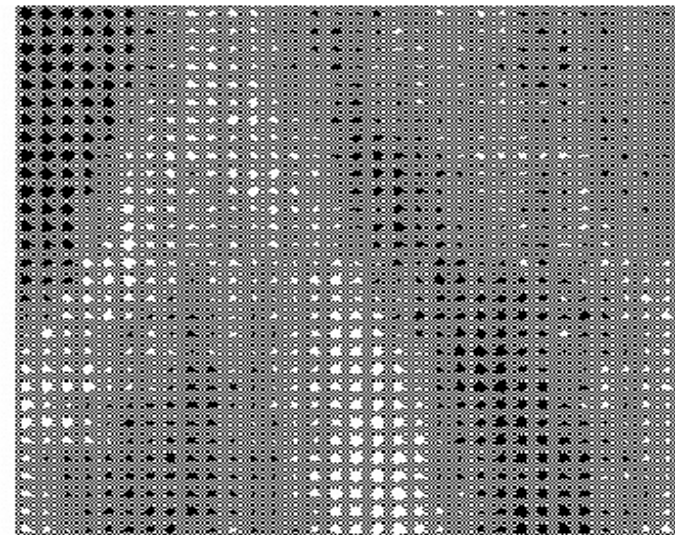


Highly non-linear decision surface

Neural Network  
trained to drive a  
car!



Weights to output units from one hidden unit



Weights of each pixel for one hidden unit

# Connectionist Models

---

Consider humans:

- Neuron switching time  $\sim .001$  second
- Number of neurons  $\sim 10^{10}$
- Connections per neuron  $\sim 10^{4-5}$
- Scene recognition time  $\sim .1$  second
- 100 inference steps doesn't seem like enough

→ much parallel computation

Properties of artificial neural nets (ANN's):

- Many neuron-like threshold switching units
- Many weighted interconnections among units
- Highly parallel, distributed process

# Expressive Capabilities of ANNs

---

## Boolean functions:

- Every boolean function can be represented by network with single hidden layer
- but might require exponential (in number of inputs) hidden units

## Continuous functions:

- Every bounded continuous function can be approximated with arbitrarily small error, by network with one hidden layer [Cybenko 1989; Hornik et al. 1989]
- Any function can be approximated to arbitrary accuracy by a network with two hidden layers [Cybenko 1988].

# **Expressive Capabilities of NNs**



# 1 hidden layer NN demo

<https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

# Prediction using Neural Networks

**Prediction** – Given neural network (hidden units and weights), use it to predict the label of a test point

**Forward Propagation** –

Start from input layer

For each subsequent layer, compute output of sigmoid unit

Sigmoid unit:

$$o(\mathbf{x}) = \sigma(w_0 + \sum_i w_i x_i)$$

1-Hidden layer,  
1 output NN:

$$o(\mathbf{x}) = \sigma \left( w_0 + \sum_h w_h \underbrace{\sigma \left( w_0^h + \sum_i w_i^h x_i \right)}_{o_h} \right)$$

# Training Neural Networks – l2 loss

$$W \leftarrow \arg \min_W E[W]$$

$$W \leftarrow \arg \min_W \sum_l (y^l - \hat{f}(x^l))^2$$

Learned neural network

Where  $\hat{f}(x^l) = o(x^l)$ , output of neural network for training point  $x^l$

Train weights of all units to minimize sum of squared errors of predicted network outputs

**Minimize using Gradient Descent**

**For Neural Networks,  
 $E[w]$  no longer convex in  $w$**

Gradient

$$\nabla E[\vec{w}] \equiv \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

i.e.,

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

# Incremental (Stochastic) Gradient Descent

---

## Batch mode Gradient Descent:

Do until satisfied

1. Compute the gradient  $\nabla E_D[\vec{w}]$

Using all training data  $D$

2.  $\vec{w} \leftarrow \vec{w} - \eta \nabla E_D[\vec{w}]$

$$E_D[\vec{w}] \equiv \frac{1}{2} \sum_{l \in D} (y^l - o^l)^2$$

---

## Incremental mode Gradient Descent:

Do until satisfied

• For each training example  $l$  in  $D$

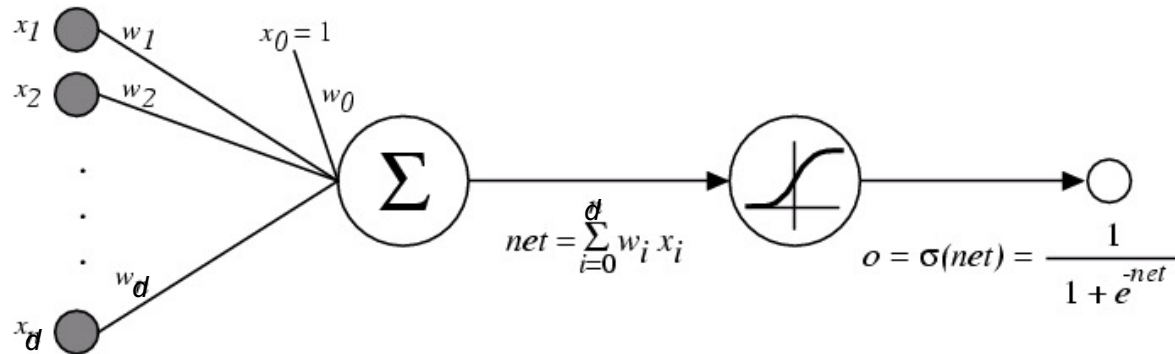
1. Compute the gradient  $\nabla E_l[\vec{w}]$

2.  $\vec{w} \leftarrow \vec{w} - \eta \nabla E_l[\vec{w}]$

$$E_l[\vec{w}] \equiv \frac{1}{2} (y^l - o^l)^2$$

*Incremental Gradient Descent* can approximate  
*Batch Gradient Descent* arbitrarily closely if  $\eta$   
made small enough

# Training Neural Networks



$\sigma(x)$  is the sigmoid function

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

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

**Differentiable**

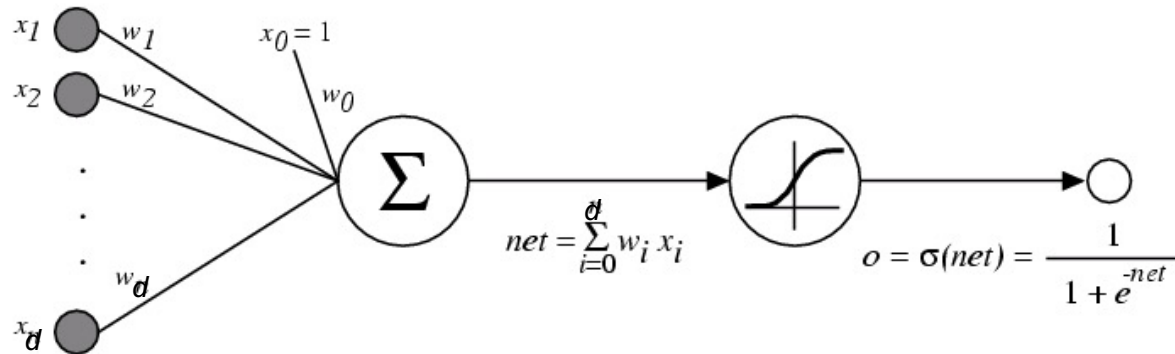
A.  $\sigma(x)(1 - \sigma(x))$

C.  $-\sigma(x)$

B.  $\sigma(x) \sigma(-x)$

D.  $\sigma(x)^2$

# Training Neural Networks



$\sigma(x)$  is the sigmoid function

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

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

We can derive gradient decent rules to train

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

# Gradient descent for training NNs

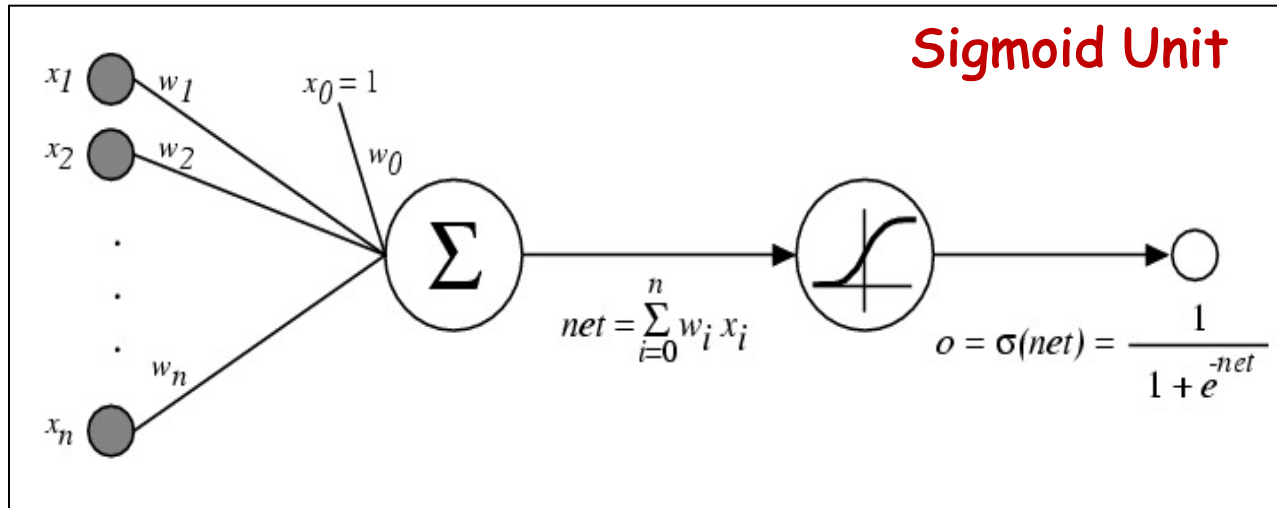
Gradient descent via Chain rule for computing gradients  
= **Back-propagation** algorithm for training NNs

Gradient of loss with respect to one weight  $w_i$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{l \in D} (y^l - o^l)^2 = \sum_l (y^l - o^l) \left( -\frac{\partial o^l}{\partial w_i} \right)$$

Stochastic gradient:

# Gradient Descent for 1 node

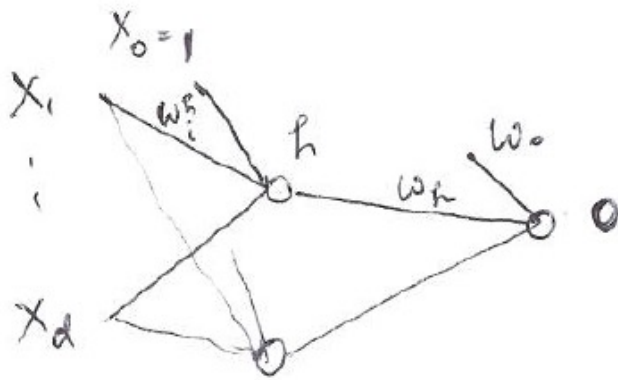


$$\frac{\partial o}{\partial w_i} = \frac{\partial o}{\partial net} \cdot \frac{\partial net}{\partial w_i} = o(1 - o)x_i$$

Gradient descent step:



# Gradient Descent for 1 hidden layer 1 output NN

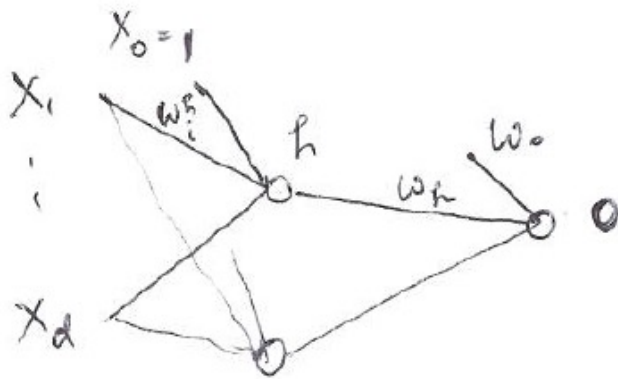


$$o = \sigma\left(w_0 + \sum_h w_h o_h\right) \equiv \sigma\left(\sum_h w_h o_h\right)$$
$$o_h = \sigma\left(w_0^h + \sum_i w_i^h x_i\right) \equiv \sigma\left(\sum_i w_i^h x_i\right)$$

Gradient of the output with respect to one **final** layer weight  $w_h$

$$\frac{\partial o}{\partial w_h} = o(1 - o)o_h$$

# Gradient Descent for 1 hidden layer 1 output NN



$$o = \sigma\left(\omega_0 + \sum_h \omega_h o_h\right) \equiv \sigma\left(\sum_h \omega_h o_h\right)$$

$$o_h = \sigma\left(\omega_0^h + \sum_i \omega_i^h x_i\right) \equiv \sigma\left(\sum_i \omega_i^h x_i\right)$$

Gradient of the output with respect to one **hidden** layer weight  $w_i^h$

$$\frac{\partial o}{\partial w_i^h} = \frac{\partial o}{\partial o_h} \cdot \frac{\partial o_h}{\partial w_i^h}$$

$$\frac{\partial o_h}{\partial w_i^h} =$$

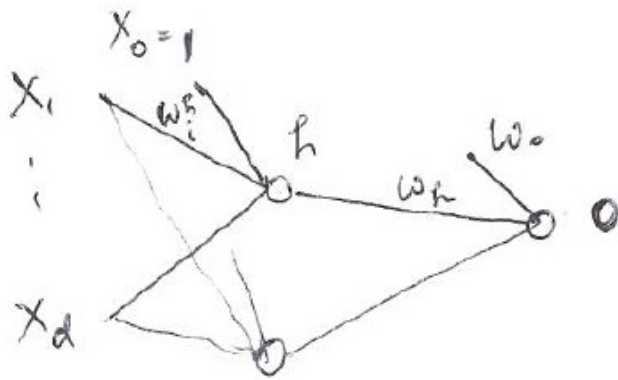
A.  $o_h(1-o_h)$

C.  $o(1-o) x_i$

B.  $o_h(1-o_h)x_i$

D.  $o(1-o)$

# Gradient Descent for 1 hidden layer 1 output NN



$$o = \sigma(w_0 + \sum_h w_h o_h) \equiv \sigma\left(\sum_h w_h o_h\right)$$

$$o_h = \sigma(w_0^h + \sum_i w_i^h x_i) \equiv \sigma\left(\sum_i w_i^h x_i\right)$$

Gradient of the output with respect to one **hidden** layer weight  $w_i^h$

$$\frac{\partial o}{\partial w_i^h} = \frac{\partial o}{\partial o_h} \cdot \frac{\partial o_h}{\partial w_i^h}$$

$$\frac{\partial o}{\partial o_h} =$$

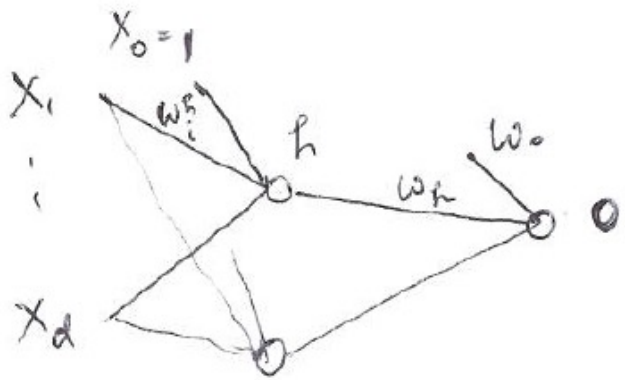
A.  $o(1-o)x_i$

C.  $o(1-o)w_h$

B.  $o(1-o)o_h$

D.  $o(1-o)$

# Gradient Descent for 1 hidden layer 1 output NN



$$o = \sigma(w_o + \sum_h w_h o_h) \equiv \sigma\left(\sum_h w_h o_h\right)$$

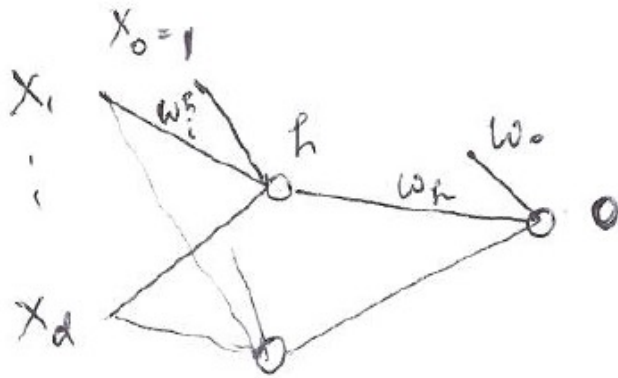
$$o_h = \sigma(w_o^h + \sum_i w_i^h x_i) \equiv \sigma\left(\sum_i w_i^h x_i\right)$$

Gradient of the output with respect to one **hidden** layer weight  $w_i^h$

$$\frac{\partial o}{\partial w_i^h} = \frac{\partial o}{\partial o_h} \cdot \frac{\partial o_h}{\partial w_i^h} = o(1 - o)w_h \cdot o_h(1 - o_h)x_i$$

$$\frac{\partial o_h}{\partial w_{i'}^h} = \frac{\partial o_h}{\partial o_h} \cdot \frac{\partial o_h}{\partial w_{i'}^h} = o_h(1 - o_h)x_{i'}$$

# Gradient Descent for 1 hidden layer 1 output NN



$$o = \sigma\left(w_0 + \sum_h w_h o_h\right) \equiv \sigma\left(\sum_h w_h o_h\right)$$

$$o_h = \sigma\left(w_0^h + \sum_i w_i^h x_i\right) \equiv \sigma\left(\sum_i w_i^h x_i\right)$$

Gradient of the output with respect to one **final** layer weight  $w_h$

$$\frac{\partial o}{\partial w_h} = o(1 - o)o_h$$

Gradient of the output with respect to one **hidden** layer weight  $w_i^h$

$$\frac{\partial o}{\partial w_i^h} = \frac{\partial o}{\partial o_h} \cdot \frac{\partial o_h}{\partial w_i^h} = o(1 - o)w_h \cdot o_h(1 - o_h)x_i$$

# Backpropagation Algorithm (MLE) using Stochastic gradient descent

1 final output unit

Initialize all weights to small random numbers.  
Until satisfied, Do

• For each training example, Do

1. Input the training example to the network  
and compute the network outputs

→ Using Forward propagation

$y$  = label of current  
training example

2.

$$\delta \leftarrow o(1 - o)(y - o)$$

3. For each hidden unit  $h$

$$\delta_h \leftarrow o_h(1 - o_h)w_h\delta$$

4. Update each network weight  $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}$$

$w_{ij}$  = wt from  $i$  to  $j$

Note: if  $i$  is input variable,  
 $o_i = x_i$

where

$$\Delta w_{i,j} = \eta \delta_j o_i$$

# Backpropagation Algorithm (MLE) using Stochastic gradient descent

Initialize all weights to small random numbers.  
Until satisfied, Do

- For each training example, Do

1. Input the training example to the network and compute the network outputs

→ Using Forward propagation

2. For each output unit  $k$

$$\delta_k \leftarrow o_k(1 - o_k)(y_k - o_k)$$

3. For each hidden unit  $h$

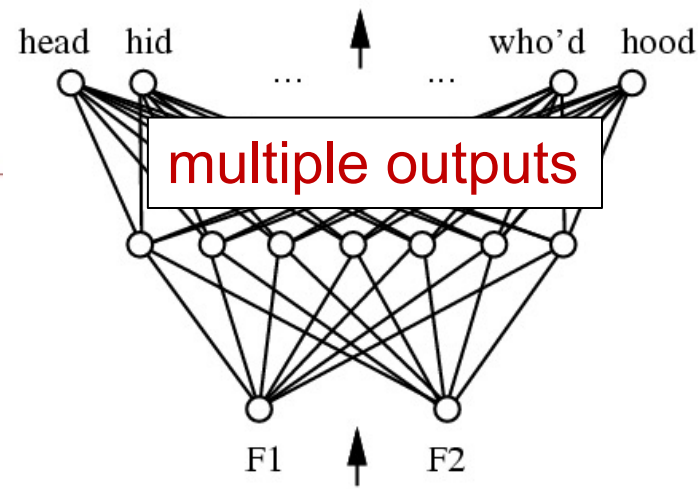
$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in \text{outputs}} w_{h,k} \delta_k$$

4. Update each network weight  $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}$$

where

$$\Delta w_{i,j} = \eta \delta_j o_i$$



$y_k$  = label of current training example for output unit  $k$

$o_{k(h)}$  = unit output (obtained by forward propagation)

$w_{ij}$  = wt from  $i$  to  $j$

Note: if  $i$  is input variable,  $o_i = x_i$

# More on Backpropagation

---

- Gradient descent over entire *network* weight vector
- Easily generalized to arbitrary directed graphs
- Will find a local, not necessarily global error minimum
  - In practice, often works well (can run multiple times)
- Minimizes error over *training* examples
  - Will it generalize well to subsequent examples?
- Training can take thousands of iterations → slow!
- Using network after training is very fast

Objective/Error no longer convex in weights