# Neural Networks to learn f: X → 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 :

**Sigmoid Unit**



$net = \sum_{i=0}^{n} w_i x_i$

$o = \sigma(net) = \frac{1}{1 + e^{-net}}$

linear combination

activation unit

Output layer, Y

Hidden layer, H

Input layer, X

① Y

½

O          X

X → [∫] → Y
   ω

①' Y

1

O          X

1 —$w_0$—\
X —$w$—) Σ → [∫] → Y
\underbrace{\hspace{3cm}}\
sigmoid unit

② Y

1

O          X

1 —$w_0$→ Σ → ∫ → +1\
X —$w$—, $w_0'$ → Σ → ∫ → -1 ) Σ → ∫ → Y\
$w'$

1

$X_1$

$X_2$

Y

① or ①'

$X_1$

$X_2$

0

1

1

Y or $P(Y|X)$

1

→ $X_2$

② 

$X_1$

0

1

0

$X_2$

1\
$X_1$\
$X_2$ → Y

# 1 hidden layer NN demo on 2D inputs

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

# 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].

# 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

error

$$\nabla E[\vec{w}] \equiv \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \cdots \frac{\partial E}{\partial w_n} \right]$$
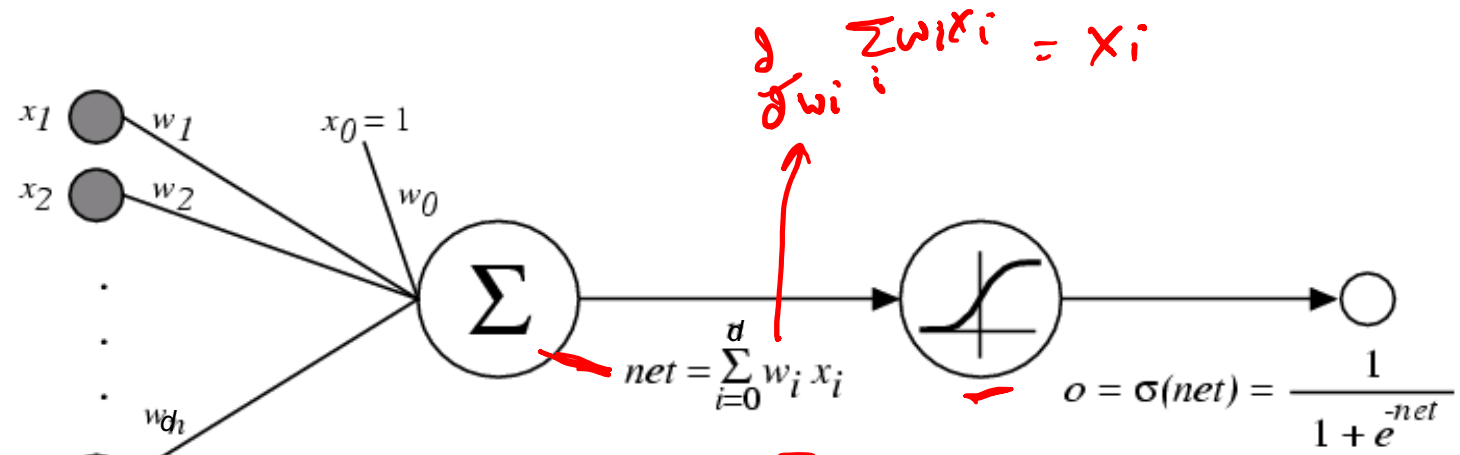
Training rule:

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

i.e.,

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

# Gradient Descent for 1 sigmoid unit



$$\frac{\partial}{\partial w_i} \sum w_i x_i = x_i$$

$$net = \sum_{i=0}^{d} w_i \, x_i$$

$$o = \sigma(net) = \frac{1}{1 + e^{-net}}$$

$$w_i \leftarrow w_i - \eta \frac{\partial E}{\partial w_i}$$

$$\frac{\partial E}{\partial o^l} \cdot \frac{\partial o^l}{\partial w_i}$$

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

$$\frac{\partial \sigma(t)}{\partial t} = \sigma(t)(1 - \sigma(t))$$

Gradient of the sigmoid function output wrt its input
$$\frac{\partial \sigma(net)}{\partial net} = \sigma(net)(1 - \sigma(net)) = o(1 - o)$$

Gradient of the sigmoid unit output wrt input weights
$$\frac{\partial o}{\partial w_i} = \frac{\partial o}{\partial net} \cdot \frac{\partial net}{\partial w_i} = o(1 - o)x_i$$

# Incremental (Stochastic) Gradient Descent

SGD

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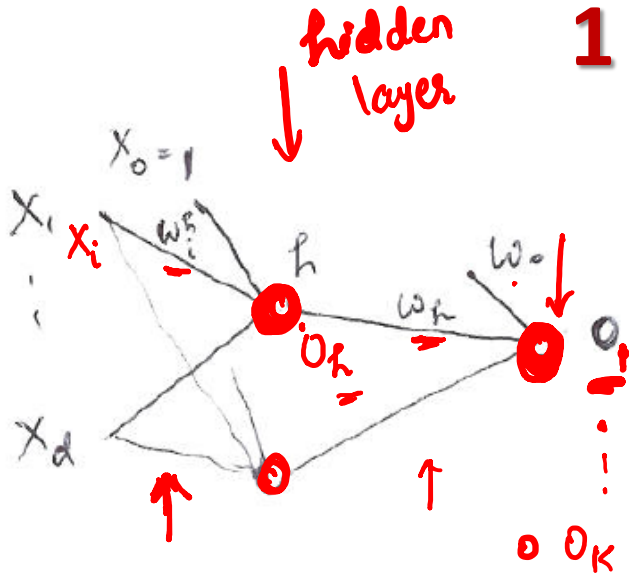
2. $\vec{w} \leftarrow \vec{w} - \eta \nabla E_l[\vec{w}]$

$$E_l[\vec{w}] \equiv \frac{1}{2}(y^l - o^l)^2$$   no sum

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

# Gradient Descent for 1 hidden layer 1 output NN

hidden layer

$X_0 = 1$

$X_1$ $X_i$ $w_i^h$ $h$ $w_o$ $O$

$O_h$ $w_h$

$X_d$

$O_K$

$$O = \sigma\left(w_o + \sum_h w_h O_h\right) \equiv \sigma\left(\sum_h w_h O_h\right)$$

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

$$\frac{\partial O}{\partial w_i^h} = \frac{\partial O}{\partial O_h} \cdot \boxed{\frac{\partial O_h}{\partial w_i^h}}$$

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

$\frac{\partial E}{\partial o^l}$

$O_h(1 - O_h) X_i$

$O(1 - O) w_h$

Gradient of the output with respect to $w_h$

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

Gradient of the output with respect to input weights $w_i^h$

$$\frac{\partial o}{\partial w_i^h} = \underbrace{o(1 - o)}_{(y-o)\,\delta} o_h(1 - o_h)w_h x_i$$

$\delta_h$

# 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 $\longrightarrow$ Using Forward propagation

  2.

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

  3. For each hidden unit $h$
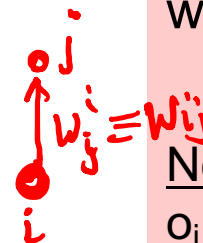
     $$\delta_h \leftarrow o_h(1 - o_h)w_h\delta$$

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

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

     where

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

$y$ = label of current training example

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

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

$o j$

$v_j^i = w_{ij}$

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

# 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 $\longrightarrow$ 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 $h$

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

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

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

  where

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

head  hid  …  ↑  …  who'd  hood

F1  ↑  F2

$y_k$ = 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, $o_i = x_i$

# HW2

$$\begin{bmatrix} O_1 \\ \vdots \\ O_k \end{bmatrix} = \begin{bmatrix} 0.1 \\ 0.2 \\ 0.6 \\ 0.05 \\ 0.05 \end{bmatrix}$$

$$\begin{matrix} O \\ O \\ O \\ O \end{matrix} \quad \begin{matrix} O_1 = P(y_1 O_1|x) \\ \vdots \\ O_k = P(y_k = O_k|x) \end{matrix}$$

➤ Classification — cross-entropy error metric

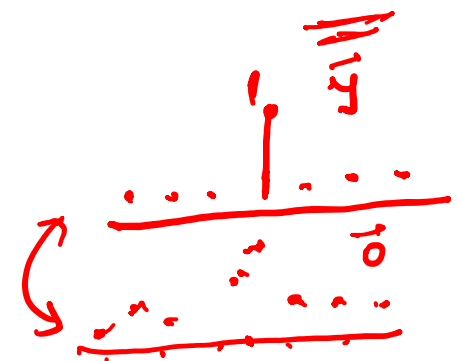$$\ell_2 : E = \frac{1}{2}(y - O)^2 \qquad \hookrightarrow E = -y \log O$$

$$-\sum_K y_k \log O_k$$

$$y_1 = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow i^{th} \text{ class}$$

$$E_y[-\log O]$$

$$-y_i \log O_i$$

$$y_k$$

$$\frac{\partial E}{\partial w} = \frac{\partial E}{\partial O} \cdot \frac{\partial O}{\partial w}$$

➤ Can implement backpropagation with matrix-vector products – uses matrix-vector calculus heavily

# Poll 1:

- $y = f(\mathbf{z})$
- $z_i = g_i(\mathbf{x})$
- $\dfrac{\partial y}{\partial \mathbf{x}} = \cdots$

A. $\dfrac{\partial y}{\partial \mathbf{z}} \dfrac{\partial \mathbf{z}}{\partial \mathbf{x}}$

B. $\dfrac{\partial \mathbf{z}^T}{\partial \mathbf{x}} \dfrac{\partial y}{\partial \mathbf{z}}$

C. $\dfrac{\partial y}{\partial \mathbf{z}} \dfrac{\partial \mathbf{z}^T}{\partial \mathbf{x}}$

D. $\dfrac{\partial y}{\partial \mathbf{z}}^T \dfrac{\partial \mathbf{z}^T}{\partial \mathbf{x}}$

E. $\left( \dfrac{\partial y}{\partial \mathbf{z}} \dfrac{\partial \mathbf{z}}{\partial \mathbf{x}} \right)^T$

F. None of the above

$$z = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \end{bmatrix} \quad 4\times1 \qquad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad 3\times1 \qquad \frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \partial y/\partial x_2 \\ \partial y/\partial x_3 \end{bmatrix} \quad 3\times1$$

$$\frac{\partial z^T}{\partial x} = \frac{\partial \overbrace{[g_1(x) \cdots g_4(x)]}^{1\times 4}}{\partial x \; 3\times 1}$$

$$= \begin{bmatrix} \frac{\partial g_1(x)}{\partial x_1} & - - & \frac{\partial g_4(x)}{\partial x_1} \\ \vdots & & \\ \frac{\partial g_1(x)}{\partial x_3} & - - - - & \frac{\partial g_4(x)}{\partial x_3} \end{bmatrix} \quad 3\times 4$$

$3\times 4$ · $4 N$ =

# 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 $\rightarrow$ slow!

- Using network after training is very fast

Objective/Error no longer convex in weights