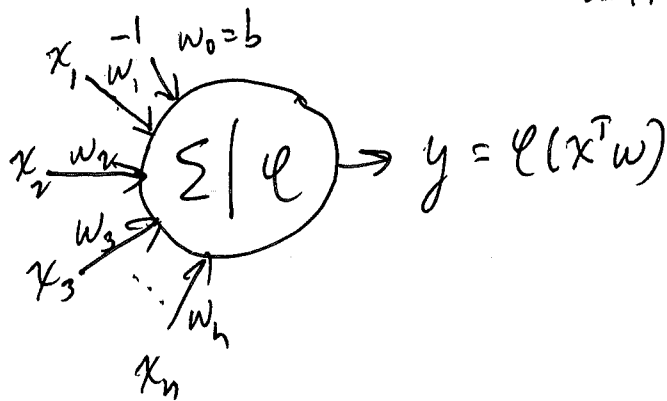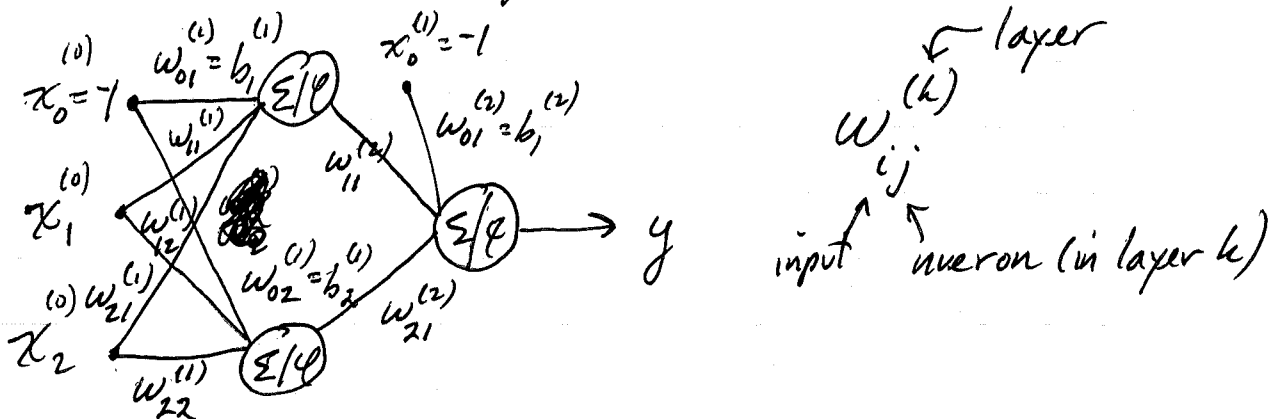Abstract neuron with bias $b$ and activation function $\ell$   [Let's try to build an artificial brain!]



$$\vec{x}\cdot\vec{w} = \vec{w}\cdot\vec{x} = \vec{x}^T\vec{w} = \vec{w}^T\vec{x} = \sum_{i=0}^{n} w_i x_i = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n - b$$

linear neuron:   $\ell(x) = x$

neural network with one hidden layer (two neurons) and two inputs

If these are linear neurons

output of : $x_1^{(1)} = \ell\left(\sum_{i=0}^{2} w_{i1}^{(1)} x_i^{(0)}\right) \overset{?}{=} \sum_{i=0}^{2} w_{i1}^{(1)} x_i^{(0)}$  (1)

neuron 1 in layer 1

"  "  " $x_2^{(1)} = \ell\left(\sum_{i=0}^{2} w_{i2}^{(1)} x_i^{(0)}\right) = \sum_{i=0}^{2} w_{i2}^{(1)} x_i^{(0)}$  (2)

" 2 " " 1

output of
neuron 1 in layer 2
or
of network!

$y = \ell\left(\sum_{i=0}^{2} w_{i1}^{(2)} x_i^{(1)}\right) = \sum_{i=0}^{2} w_{i1}^{(2)} x_i^{(1)}$  (3)

Using (1) and (2) in (3)

$\cancel{y = \left(\sum_{i=0}^{2} w_{i1}^{(2)} \sum_{i=0}^{2} w_{i1}^{(1)} x_i^{(0)}\right)}$

$y = w_{01}^{(2)} x_0^{(1)} + w_{11}^{(2)} x_1^{(1)} + w_{21}^{(2)} x_2^{(1)}$

$= w_{01}^{(2)} x_0^{(1)} + w_{11}^{(2)} \sum_{i=0}^{2} w_{i1}^{(1)} x_i^{(0)} + w_{21}^{(2)} \sum_{i=0}^{2} w_{i2}^{(1)} x_i^{(0)}$
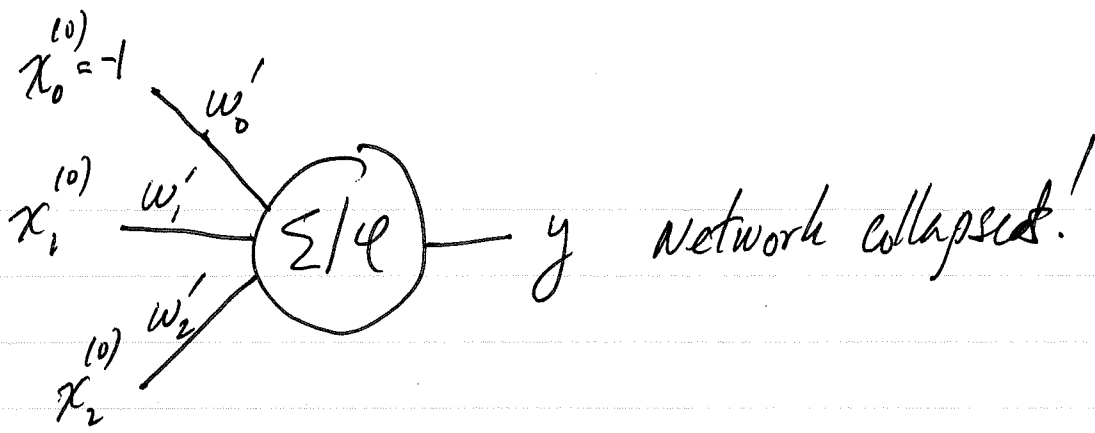
$= \underbrace{\sum_{i=0}^{2} \left(w_{11}^{(2)} w_{i1}^{(1)} + w_{21}^{(2)} w_{i2}^{(1)}\right)}_{w_i'} x_i^{(0)} + w_{01}^{(2)} x_0^{(1)}$

$= \underbrace{\left(w_0' x_0^{(0)} + w_{01}^{(2)} x_0^{(1)}\right)}_{\text{constant!}} + w_1' x_1^{(0)} + w_2' x_2^{(0)}$   Linear!

$w_0' = b'$

What is a new simplified
network representation?

$x_0^{(0)} = -1$    $w_0'$

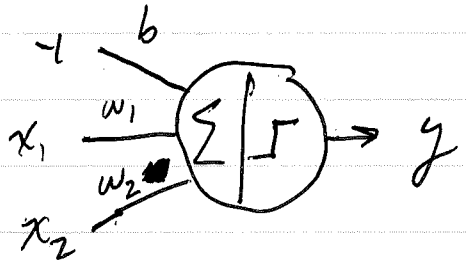$x_1^{(0)}$    $w_1'$    $\Sigma/\varphi$    $y$    Network collapsed!

$w_2'$

$x_2^{(0)}$

$\varphi$ needs to be non-linear!

Perceptron    $\varphi(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$    $\varphi(x)$

Assume

$-1$    $b$

$x_1$    $w_1$    $\Sigma/\Gamma$    $y$

$w_2$

$x_2$

which of the following logic gates can this implement (if any)?

AND

| $x_1$ | $x_2$ | $y = x_1 \wedge x_2$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

OR

| $x_1$ | $x_2$ | $y = x_1 \vee x_2$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

XOR

| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

✓ AND?

$x_2$ (0,1)     $x_1 + x_2 = 1.5$

(1,1)

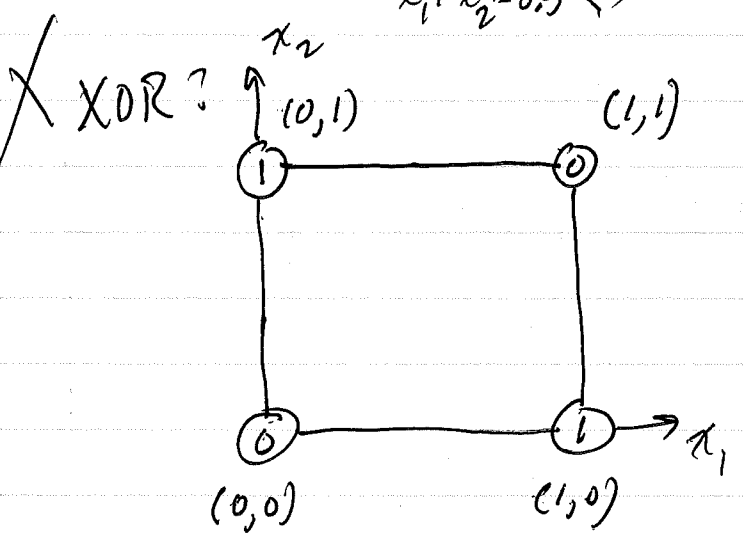$x_1 + x_2 = 1.5$

(0,0)     $x_1 + x_2 < 1.5$     (1,0)     $x_1$

let $w_1 = w_2 = 1$ and $b = 1.5$

$$y = \varphi(x_1 + x_2 - 1.5) = \begin{cases} 0, \text{ if } x_1 + x_2 < 1.5 \\ 1, \text{ if } x_1 + x_2 \geq 1.5 \end{cases}$$

1

0

✓ OR?     (0,1)     (1,1)

$x_2$

$x_1 + x_2 \geq 0.5$

let $w_1 = w_2 = 1$ and $b = 0.5$

$x_1$

(0,0)     (1,0)

$x_1 + x_2 < 0.5$

$x_1 + x_2 = 0.5$
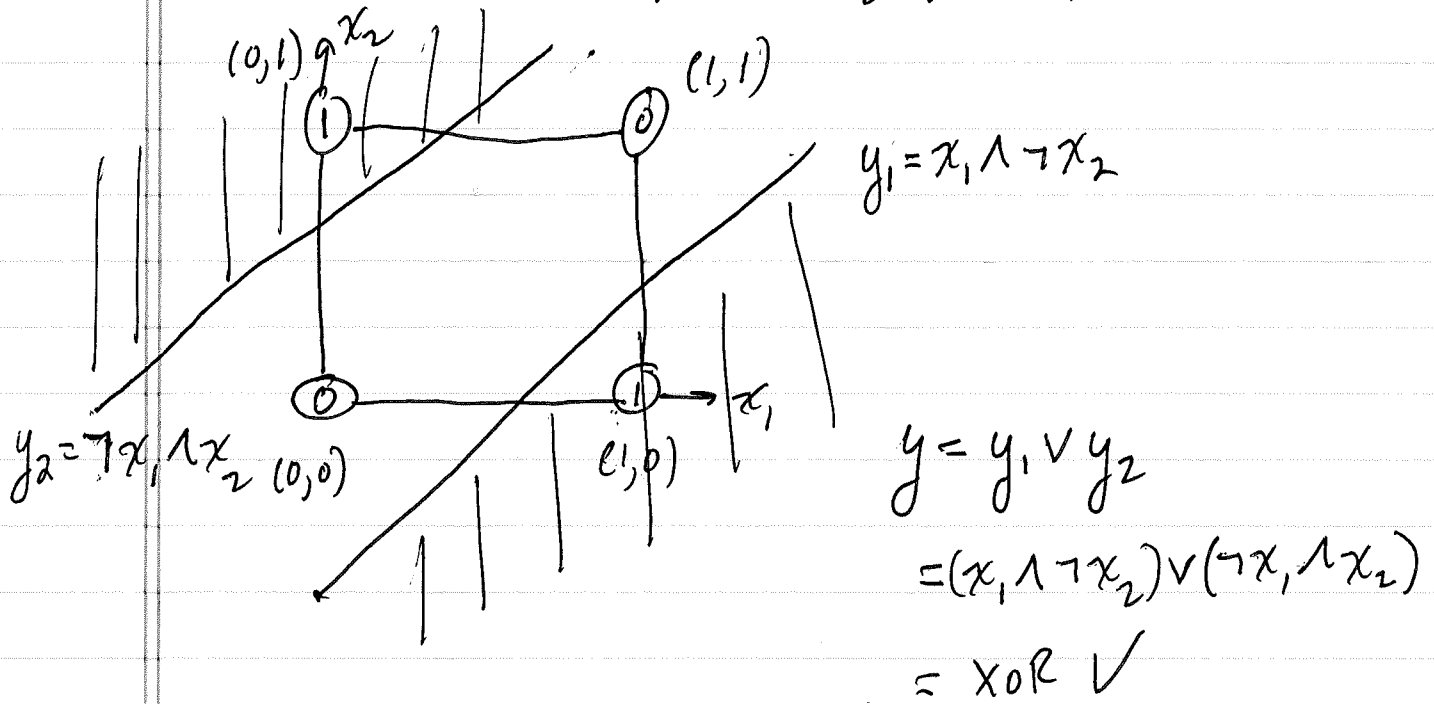
$$y = \varphi(x_1 + x_2 - 0.5) = \begin{cases} 0, \text{ if } x_1 + x_2 < 0.5 \\ 1, \text{ if } x_1 + x_2 \geq 0.5 \end{cases}$$

✗ XOR?     $x_2$

(0,1) ①     (1,1) ⓪

(0,0) ⓪     (1,0) ①     $x_1$

No line can separate the 0 and 1 output classes!

what about our 3 neuron model with 5
linear neurons replaced by perceptrons?



$(0,1)$ $\quad x_2$

$(1,1)$

$y_1 = x_1 \wedge \neg x_2$

$y_2 = \neg x_1 \wedge x_2$ $\quad (0,0)$

$(1,0)$ $\quad x_1$

$y = y_1 \vee y_2$

$$= (x_1 \wedge \neg x_2) \vee (\neg x_1 \wedge x_2)$$

$$= XOR \vee$$

Non-linear neurons and multiple/hidden
layered network yields complex decision
boundaries!

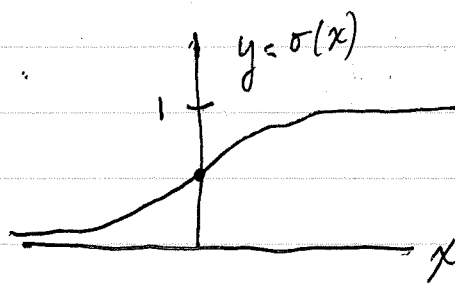Universal Function Approximator!



$x$

$f(x)$

note: $O(m^2)$ fr SVMs

So are SVMs but NN training $O(m)$ m samples

What are some other non-linear neurons?

## Sigmoid

$$\varphi(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$$



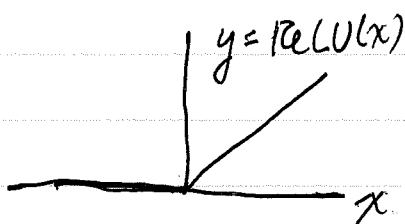What does this represent: $y = \sigma(x^T w)$?

Logistic Regression! $\Rightarrow$ 1 Sigmoid neuron

Early versions of NNs used sigmoid neurons, but because they <u>saturate</u> lead to inability to learn! Known as "vanishing gradient problem" or "barren plateaus".

How to solve this problem? ReLU(x)!

Rectified Linear Unit (ReLU)

$$\varphi(x) = ReLU(x) = \max\{0, 0\} = \begin{cases} 0, & \text{if } x < 0 \\ x, & \text{if } x \geq 0 \end{cases}$$

For multi-class output that
takes a vector $\vec{x} = (x_1, \dots, x_K)$
as input, ● use softmax

$$\varphi(\vec{x}) = \sigma(\vec{x}) = \left( \underbrace{\frac{e^{x_1}}{\sum_{k=1}^{K} e^{x_k}}}_{p(y=1|\vec{x})}, \dots, \underbrace{\frac{e^{x_K}}{\sum_{k=1}^{K} e^{x_k}}}_{p(y=K|\vec{x})} \right)$$

$$\sigma(\vec{x})_i = p(y=k|\vec{x})$$

$$\sum_{k=1}^{K} p(y=k|\vec{x}) = \sum_{k=1}^{K} \frac{e^{x_k}}{\sum_{k'=1}^{K} e^{x_{k'}}} = \frac{\sum_{k=1}^{K} e^{x_k}}{\sum_{k'=1}^{K} e^{x_{k'}}} = 1$$

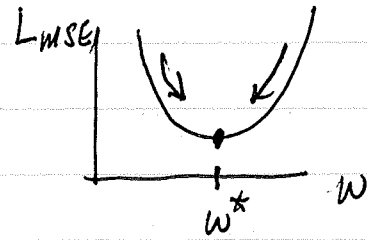where $\log \sigma(\vec{x})_i = x_i - \log \sum_{k=1}^{K} e^{x_k}$

This avoids vanishing gradient problem!

How do we train a NN for
an arbitrary dataset?

How do we fit a linear model
with multiple parameters?

Find the weights $\vec{w}$ that minimize the (MSE) mean square error <u>loss</u>

$$L_{MSE} = \frac{1}{M} \sum_{i=1}^{M} \left( y_i - f_w(x^i) \right)^2$$



where $f_w(x) = w_0 + w_1 x_1 + w_2 x_2 + \cdots + w_n x_n$

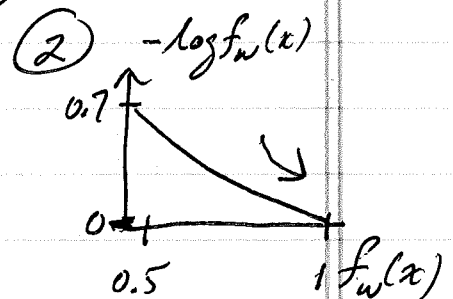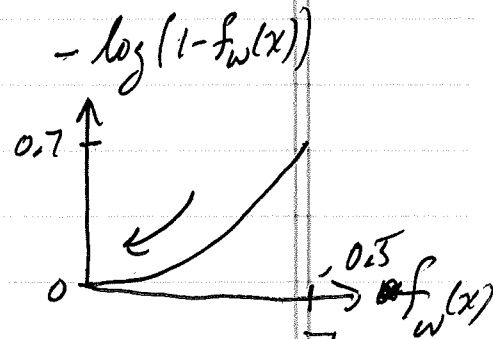Do the same for NNs! MSE is used when output of NN is real number.

What is a good loss function for a classifier (e.g. binary classifier) where output is 0 or 1?

Binary Cross Entropy:

$$L_{BCE} = \frac{-1}{M} \sum_{i=1}^{M} \left[ y_i \log f_w(x^i) + (1-y_i) \log\left(1 - f_w(x^i)\right) \right]$$



if $y_i = 1$   ② $= 0$   ① $= -\log f_w(x^i)$

if $y_i = 0$   ① $= 0$   ② $= -\log\left(1 - f_w(x^i)\right)$

How do we minimize loss functions for a neural network?

Method of Steepest Descent or Gradient Descent



Taylor Expansion (Linear Approximation):

$$\Delta L = L(w^0 + \eta v) - L(w^0) = \sum_{i=1}^{n} \frac{\partial L}{\partial w_k}(w^0) \eta v^k + O(\eta^2)$$

change weights from $w^0$ in $v$ ~~direc~~ (unit vector) direction with step size $\eta$ (which we assume is small).

$$\Delta L = \eta \vec{\nabla}L \cdot \vec{v} + O(\eta^2)$$

↑ projection of $\vec{\nabla}L$ onto $\vec{v}$

Cauchy inequality

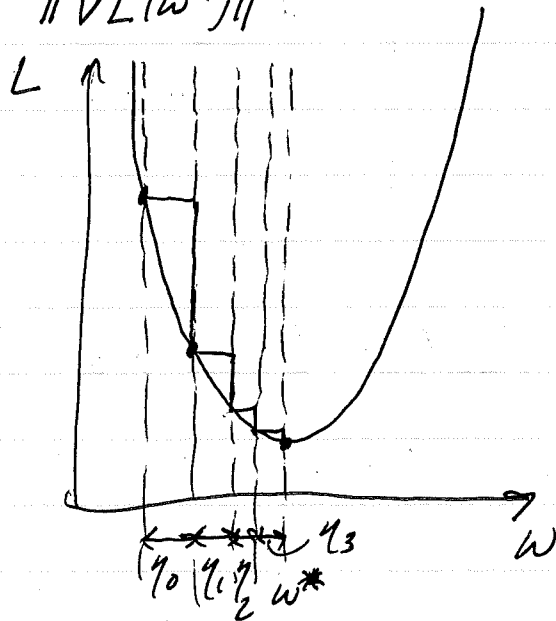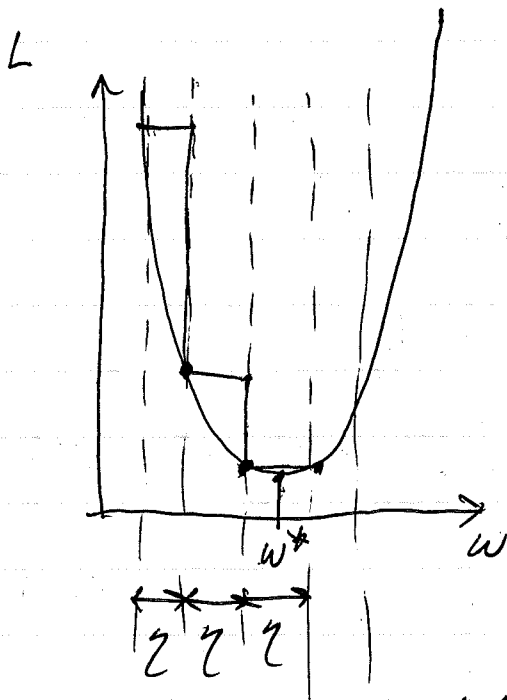$$-\|\vec{\nabla}L(w^0)\| \|\vec{v}\| \leq \vec{\nabla}L \cdot \vec{v} \leq \|\vec{\nabla}L(w^0)\| \|\vec{v}\|$$

$\vec{\nabla}L$ & $\vec{v}$ anti-parallel     $\vec{\nabla}L$ & $\vec{v}$ parallel

Maximum drop in L occurs for unit length
$\vec{v}$ ($\Rightarrow \|\vec{v}\|=1$) when it's anti-parallel to $\vec{\nabla}L$

$$\vec{v} = - \frac{\vec{\nabla}L(w^0)}{\|\vec{\nabla}L(w^0)\|}$$

$\therefore$ Update weights as follows

$$w^{n+1} = w^n - \eta \frac{\vec{\nabla}L(w^n)}{\|\vec{\nabla}L(w^n)\|}$$



Problem: Overshoots!

Fixed: Learning rate $\propto$ gradient

$$\text{let } \eta_n = \delta \|\vec{\nabla}L(w^n)\|$$

$$w^{n+1} = w^n - \eta_n \frac{\vec{\nabla}L(w^n)}{\|\vec{\nabla}L(w^n)\|} = w^n - \delta\|\vec{\nabla}L(w^n)\| \frac{\vec{\nabla}L(w^n)}{\|\vec{\nabla}L(w^n)\|}$$

$$w^{n+1} = w^n - \delta\vec{\nabla}L(w^n)$$

How calc $\nabla L$ w/ $f_w(x^i)$

representing a NN? Back Propagation