

→ Perceptron is a linear classifier

→ The  $\vec{w}_s$  will be the weights (as earlier).  
 $w_0$  is the intercept  $w_0 = 1$  initially.

A perceptron calculates ② quantities  $\begin{cases} \text{① a weighted sum of the input features} \\ \text{② a threshold on this sum by the ① function.} \end{cases}$

The perceptron is a simple artificial model of human neurons

weights  $\rightarrow$  synapses

threshold  $\rightarrow$  neuron firing

Perceptron operation:  $\hat{f}(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } \theta_0 x_0 + \dots + \theta_n x_n > 0 \\ 0 & \text{otherwise (can also be defined as -1)} \end{cases}$

$\vec{\theta} \rightarrow$  weight vector  $\in \mathbb{R}^{n+1}$

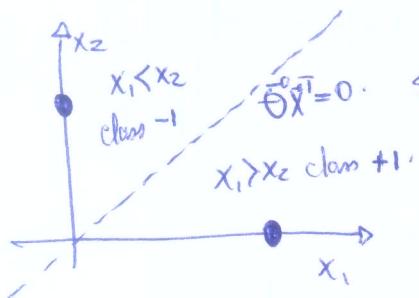
$\vec{x} \rightarrow$  feature vector  $\in \mathbb{R}^{n+1}$

$$\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n = \underbrace{\vec{\theta}^T \vec{x}}_{\text{vector inner product}}$$

The perceptron represents a hyperplane decision surface in  $n$ -dimensional space  
 } 2d line  
 } 3d plane

The equation of the hyperplane is  $\vec{\theta}^T \vec{x} = 0$   
 This is the equation for points in  $x$ -space that are on the boundary.

Example



Suppose  $\vec{\theta} = (\theta_0 = 1, \theta_1 = 1, \theta_2 = -1)$ .

$$\vec{\theta}^T \vec{x} \Rightarrow 1x_0 + 1x_1 + (-1)x_2 = 0$$

$$x_1 - x_2 = 0$$

$x_1 = x_2$  decision boundary equation

- \* A data set is ~~not~~ separable by a learner if there is some instance of that learner that correctly predicts all data points (37)

- \* The learner is linearly separable if:

- Can separate the two classes using a straight line if  $x \in \mathbb{R}^2$  or with a hyperplane if  $x \in \mathbb{R}^n$ .

### Class overlap

~~separable~~ } Common in practice  
some observation values possible for both classes (benign/malignant cells look similar).

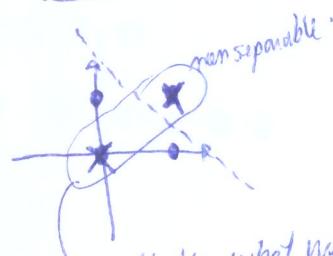
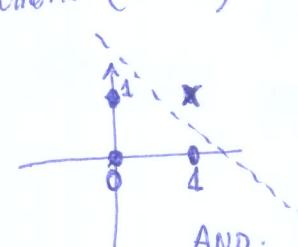
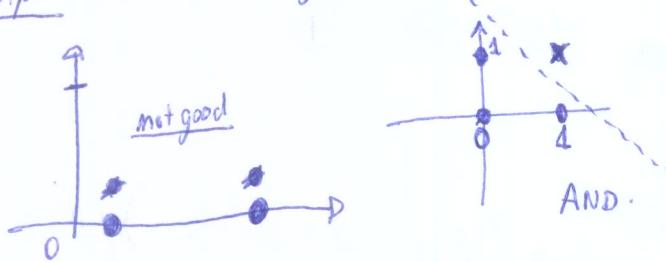
what to do? } - Increase number of features.  
- Increase model complexity.

### Representational power of perceptions

so what mappings a perception can represent perfectly?

well, it is a linear classifier, so it can represent any mapping that is linearly separable

Example: some boolean functions (AND) but no XOR



↑ that's what we would need to learn.

### Remember

### Effects of dimensionality

① Data are increasingly separable in high dimension - Is this a good thing?

Good - separation is easier in higher dimensions (for fixed  $M$ )

Good - increase nbr of features, and even a linear classifier would do the trick!

Combe examples.

Bad - training vs test error; overfitting; bias vs variance

Bad - increasingly complex decision boundaries can eventually get all training data right but it doesn't mean necessarily good for testing data  $\Rightarrow$  lacks generalization

How to update weights?

for each data point  $x^{(i)}$ :

$$f(x^{(i)}) \leftarrow T(\tilde{w} * x^{(i)})$$

$$\tilde{w} \leftarrow w + \alpha (f(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \quad \text{gradient-like step}$$

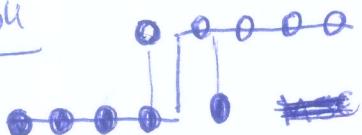
↳ if  $f(x^{(i)})$  correctly predicted: no update.  
otherwise: update weights.

Another way of thresholding would be using smooth functions



} if we are far from decision boundary,  $|f(\cdot)|$  is large, small error.  
} nearby the boundary:  $|f(\cdot)|$  near  $\frac{1}{2}$ , large error.

Example



$$J(\theta) = \frac{\text{sum of errors}}{10} \quad \text{normal thresholding.}$$

while



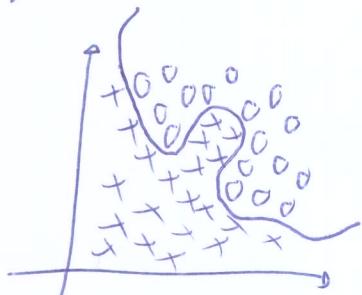
$$J(\theta) = \frac{(0^2 + 1^2 + 2^2 + 25^2 + \dots)}{10}.$$

summary - linear classifier  $\Leftrightarrow$  perceptron

- visualizing the decision boundary: useful but not always possible
- measuring quality: sum of squared errors (SSE).
- Learning the weights of a linear classifier reduces to an optimization problem. For SSE we can do gradient descent but there are others

- with old idea
- out of favor for sometime but in the hyp again for several recent results such as Deep Belief Networks (~~deep learning~~) and results such as image into
- why do we need yet another learning algorithm?

Consider a problem



→ we could use LogReg such as

$$\hat{f}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \dots)$$

poly terms.

→ normally it works OK for a few features (2, 3) but not as good for many features. (how to derive polynomials thru?).

if we have 100 features and want to just include 2nd order polynomials, we would have

$$\hat{f}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_{100} x_{100} + \theta_{101} x_1^2 + \theta_{102} x_1 x_2, \dots, \theta_k x_1 x_{100} + \theta_{k+1} x_2 x_3 + \dots) \approx 5K \text{ features. } \Theta(n^2) \approx \frac{n^2}{2}$$

≈ If we want 3rd order polys, we would have  $\Theta(n^3)$ .

For 100 features  $\approx 170K$  features!

≈ In practice, we have problems with thousands or millions of features in the original ~~feature~~ feature space.

## Neurons and the brain

### Origins

{ algorithms that try to mimic the brain.  
widely used in the 80's and early 90's. Popularity went down in late 90's.  
Recent hype due to state-of-the-art results in many applications.  
the reason is that now we have enough computational power to design and use  
large-scale neural ~~netw~~ networks

expensive to  
run

David

## The "one learning algorithm" hypothesis

→ some people believe the brain does all of its amazing things based on one learning algorithm.



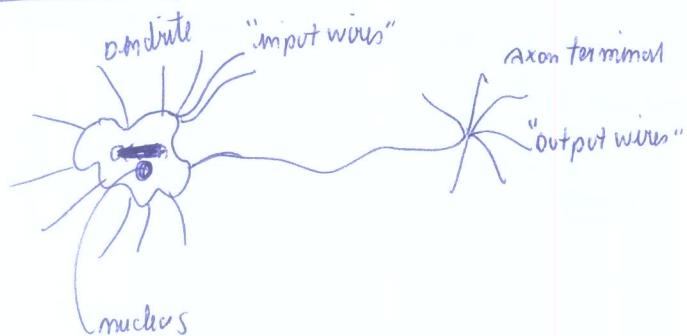
Auditory cortex learns to ~~see~~ hear/listen  
but researchers have shown that cutting  
the connection of Auditory cortex to the eye  
and rewiring it to the ear leads to a learning  
process in which the auditory cortex learns  
to see!

### neuro-rewiring experiments

Examples :- Bramport experiment: Camera for low-res grayscale image on top of the head and a wire to connect each pixel to our tongue (low voltage → dark pixel and we can "see" through our tongue). high voltage → bright pixel.

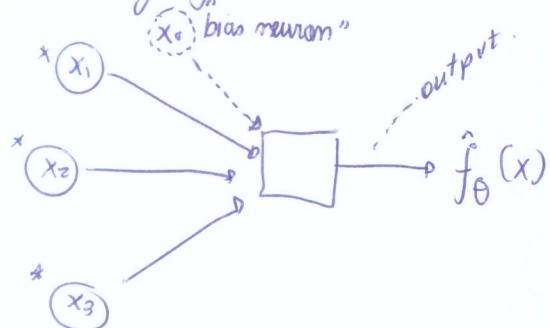
- Human echolocation (sonar): snap fingers. Boy that had his eyeballs removed learned to navigate using this "rewiring" technique.
- Implanting a 3rd eye in a frog will lead him to actually learn how to use it.

## Neural Networks - Representations



Computationally speaking, we can think of a neuron as a computational unit that receives a set of inputs, performs some computation and outputs some signals through the axons to other neurons.

We are going to model this as a logistic unit



\*  $\Rightarrow$  inputs.

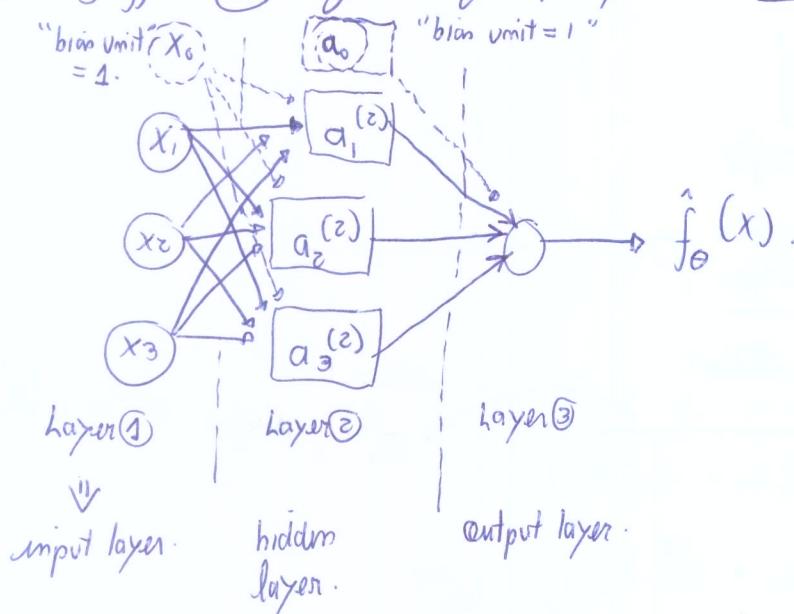
$$\hat{f}_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$\vec{x} = \begin{bmatrix} x_0 \\ \vdots \\ x_m \end{bmatrix} \quad \vec{\theta} = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_m \end{bmatrix}$$

$x_0 \Rightarrow$  bias neuron always equal to 1.

In NN terminology the sigmoid/logistic activation function and the  $\vec{\theta}$  parameters are known as weights.  $\vec{\theta}$  represents

Generalizing, a NN is just a group of several neurons.



because it is not  $\vec{x}$  or  $\vec{y}$  (input/output)  
therefore the term "hidden".

→ We may have NN with several hidden layers.

$a_i^{(j)}$  = "activation" of unit  $i$  in layer  $j$

$\Theta^{(j)}$  = matrix of weights controlling function mapping from layer  $(j)$  to layer  $(j+1)$ .

Here's the computation:

$$a_1^{(2)} = g\left(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3\right)$$

$$a_2^{(2)} = g\left(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3\right).$$

$$a_3^{(2)} = g\left(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3\right).$$

$$\text{if } f_\theta(x) = a_1^{(3)} = g\left(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)}\right).$$

$\Theta^{(1)}$  is a matrix mapping layer  $(1)$  to  $(2)$ , so  $\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$ .

→ If network has  $s_j$  units in layer  $(j)$ , and  $s_{j+1}$  in layer  $(j+1)$ , then  $\Theta^{(j)}$  will have dimensions  $(s_{j+1}) \times (s_j + 1)$ .

→ with this, we have  $\hat{f}_\theta(x)$   
our classifi.

model representation, intuition and calculations

forward propagation: vectorized implementation

$$a_1^{(2)} = g\left(\underbrace{\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3}_{z_1^{(2)}}\right)$$

$$a_2^{(2)} = g\left(\underbrace{\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3}_{z_2^{(2)}}\right)$$

$$a_3^{(2)} = g\left(\underbrace{\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3}_{z_3^{(2)}}\right)$$

let's call  $z_i^{(2)}$ , therefore  $a_i^{(2)} = g(z_i^{(2)})$ .

$$\boxed{\Theta^{(1)} \cdot X} \text{ or } \boxed{\Theta^{(1)} \cdot \tilde{X}}$$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_3 \end{bmatrix} ; Z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} ; \quad \cancel{\text{...}}$$

$$z^{(2)} = \vec{\theta}^{(1)} \vec{x}$$

$$a^{(2)} = g(z^{(2)}). \quad \begin{matrix} \text{in } \mathbb{R}^3 \\ \text{out } \mathbb{R}^3 \end{matrix}$$

(g) function applies Sigmoid function on each  $z_i^{(2)}$  element-wise.  
Given all elements in  $z^{(2)}$  (e.g. in  $\mathbb{R}^3$ ), we apply (g) on each individually.

Let's define  $a^{(1)} = \vec{x}$ , therefore  $z^{(2)} = \vec{\theta}^{(1)} \cdot \vec{x}$  now becomes

$$z^{(2)} = \vec{\theta}^{(1)} \cdot a^{(1)}$$

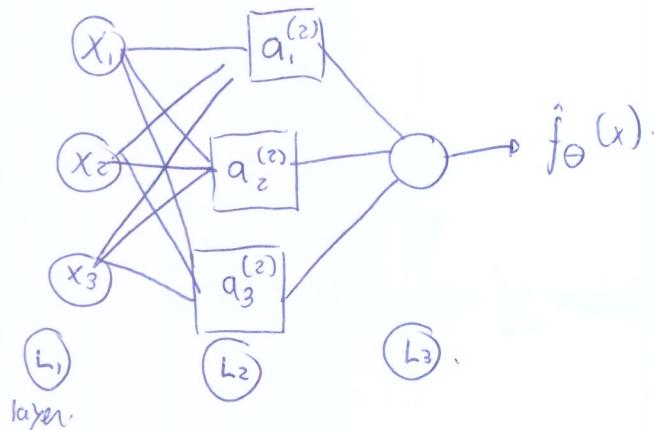
Now we [add]  $a_0^{(2)} = 1$  so that  $a^{(2)} \in \mathbb{R}^4$   
bias unit and

$$z^{(3)} = \vec{\theta}^{(2)} \cdot a^{(2)}$$

$$\hat{f}_{\theta}(x) = a^{(3)} = g(z^{(3)})$$

This algorithm is known as forward propagation.

Neural Network Learning its own features



\* Suppose we don't see the left part of it.

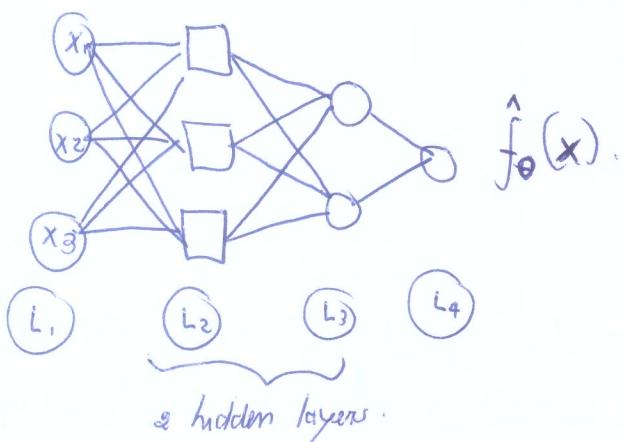
if we do that, we end up having a logistic Regression-like situation

$$\hat{f}_{\theta}(x) = g\left(\Theta_{00}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{22}^{(2)} a_2^{(2)} + \Theta_{33}^{(2)} a_3^{(2)}\right).$$

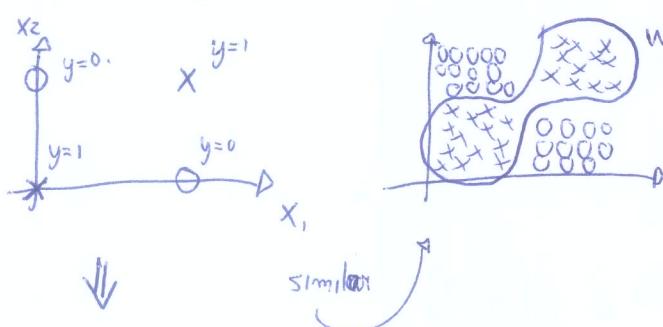
The cool thing here is that the features  $(a_i)$  are learned from the input.

Observation: we can have other NN diagrams/architectures

↳ how ~~NN~~ units connect to each other.



Let's return to ~~our~~ logical operators example in which  $(x_1)$  and  $(x_2)$  are binary (0 or 1).



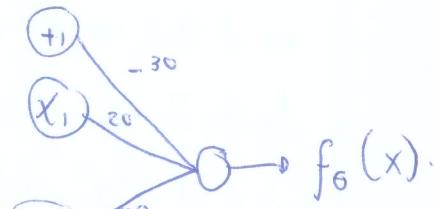
We want to learn a non-linear separator/classifier

$$\begin{cases} y = x_1 \text{ XOR } x_2 \\ y = \text{XNOR } x_1 \\ \text{Not}(x_1 \text{ XOR } x_2) \end{cases}$$

Let's start with a simple example for logical AND.

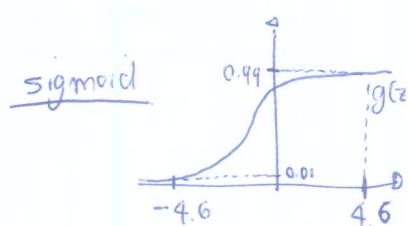
$$x_1, x_2 \in \{0, 1\}$$

$$y = x_1 \text{ AND } x_2.$$



$$\hat{f}_\theta(x) = g(-30x_0 + 20x_1 + 20x_2)$$

$$\theta_{10}^{(1)} \quad \theta_{11}^{(1)} \quad \theta_{12}^{(1)}$$

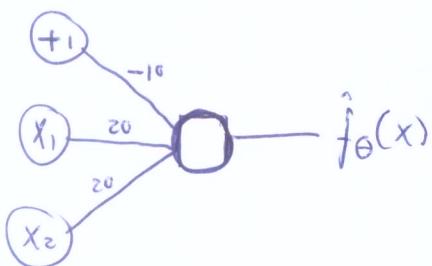


$x_1$	$x_2$	$\hat{f}_\theta(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

$$\hat{f}_\theta(x) \approx x_1 \text{ AND } x_2$$

## Example for OR function

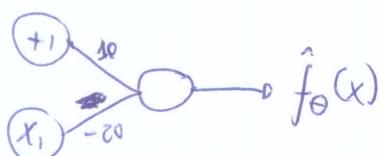
45



$x_1$	$x_2$	$\hat{f}_\theta(x)$
0	0	$g(-10) \approx 0$
0	1	$g(+10) \approx 1$
1	0	$g(+10) \approx 1$
1	1	$g(40) \approx 1$

$$\hat{f}_\theta(x) = g(-10x_0 + 20x_1 + 20x_2).$$

## Example for NOT $x_1$



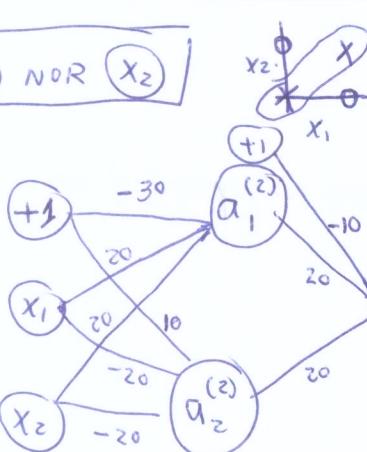
$$\hat{f}_\theta(x) = g(10 - 20x_1)$$

$x_1$	$f_\theta(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

Exercise: how to compute

$$\rightsquigarrow (\text{Not}(x_1)) \text{AND} (\text{Not}(x_2)) ?$$

$(x_1) \text{ NOR } (x_2)$



$x_1$	$x_2$	$a_1^{(2)}$	$a_2^{(2)}$	$\hat{f}_\theta(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

$\left\{ \begin{array}{l} a_1^{(2)} \text{ is using } x_1 \text{ AND } x_2 \text{ network} \\ a_2^{(2)} \text{ is using } \neg x_1 \text{ and } \neg x_2 \text{ or } (\text{not } x_1) \text{ AND } (\text{not } x_2) \end{array} \right.$

$\left\{ \begin{array}{l} a_1^{(2)} \text{ is using } x_1 \text{ AND } x_2 \text{ network} \\ a_2^{(2)} \text{ is using } \neg x_1 \text{ and } \neg x_2 \text{ or } (\text{not } x_1) \text{ AND } (\text{not } x_2) \end{array} \right.$

\* Each layer in the sequence can learn more complex functions ~~than~~ on top of its predecessors output layers.

match

Multiclass Classification

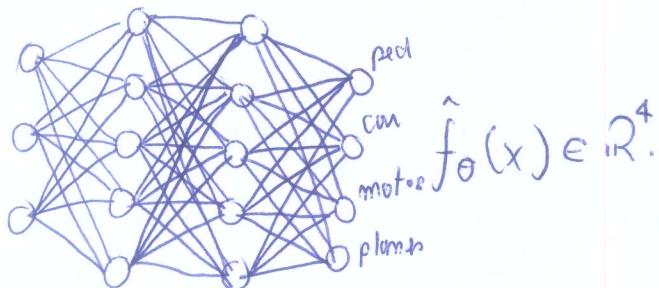
Basically, it is an extension of the one vs all method.

Pedestrian

Car

motorcycles

Planes



We want  $\hat{f}_\theta(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$  when pedestrian,  $\hat{f}_\theta(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$  when car,  $\hat{f}_\theta(x) = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$  when motorcycles and  $\hat{f}_\theta(x) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$  when plane.

In the training set,  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots, (x^{(m)}, y^{(m)})$

$y^{(i)}$  is going to be one of  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

ped      car      motor      plane

As previously, we have represented  $y \in \{1, 2, 3, 4\}$  (e.g., k-NN). Now we represent it in a binary form. So, for 4 classes,

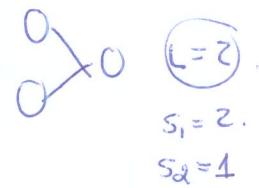
$$\hat{f}_\theta(x) \approx y^{(i)} \in \mathbb{R}^4$$

$$x$$

## Cost function for Neural Networks

Let's start with the classification problem.

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}.$$



$(L)$  will be the number of layers in the network.

$(S_l)$  = number of units (excluding the "bias" unit) in layer  $(l)$ .

In the binary classification,  $y=0$  or  $y=1$ , therefore we will have  $(1)$  output unit. For a multiclass classification problem ( $K$  classes),  $y \in \mathbb{R}^K$ , therefore we will have  $(K)$  output units.

Let's now define a cost function. Recall that for logistic regression, we had:

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log \hat{f}_\theta(x^{(i)}) + (1-y^{(i)}) \cdot \log (1-\hat{f}_\theta(x^{(i)})) \right] + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2}_{\text{Regularization}} \quad \theta_0 = 1$$

for a neural network, instead of having just one logistic regression output unit, we will have  $(K)$  of them.

$\hat{f}_\theta(x) \in \mathbb{R}^K \quad (\hat{f}_\theta(x))_i = i^{\text{th}}$  output.

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log (\hat{f}_\theta(x^{(i)}))_k + (1-y_k^{(i)}) \cdot \log (1-\hat{f}_\theta(x^{(i)}))_k \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

Wow! How can we optimize such cost function to actually find our weights or parameters?

## Back propagation Algorithm

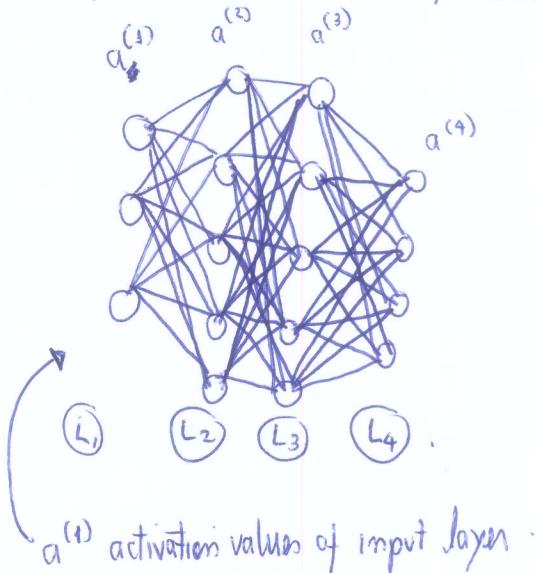
Given  $\theta$  and  $J(\theta)$ , we want to  $\min_{\theta} J(\theta)$ , for that we need to compute:

$$\textcircled{1} \quad J(\theta)$$

$$\textcircled{2} \quad \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) ; \text{ recall } \theta_{ij}^{(l)} \in \mathbb{R}.$$

Suppose we only have one training example  $(x, y)$ . Using the forward propagation, we have

$$\left\{ \begin{array}{l} a^{(1)} = x \\ z^{(2)} = \Theta^{(2)} \cdot a^{(1)} \\ a^{(2)} = g(z^{(2)}) \text{ add } a_0^{(2)} \\ z^{(3)} = \Theta^{(3)} \cdot a^{(2)} \\ a^{(3)} = g(z^{(3)}) \cdot \text{add } a_0^{(3)} \\ z^{(4)} = \Theta^{(4)} \cdot a^{(3)} \\ a^{(4)} = f_{\theta}(x) = g(z^{(4)}). \end{array} \right.$$



### Gradient computation: Back propagation algorithm

Intuition  $\delta_j^{(l)} = \text{"error" of node } j \text{ in layer } l$ .

$a_j^{(l)}$   $\Rightarrow$  activation of  $j$ th unit in layer  $l$ .

Concretely, if we have the NN above, for each output unit (layer  $L=4$ ),

$$\delta_j^{(4)} = \begin{bmatrix} a_j^{(4)} - y_j \\ \vdots \\ \delta_j^{(4)} \end{bmatrix} \quad \left\{ \begin{array}{l} \text{if we vectorize it, we have } \delta^{(4)} = \vec{a}^{(4)} - \vec{y} \\ \vec{y} \in \mathbb{R}^L \end{array} \right. \quad \text{where } L \text{ is the number of output units in the NN.}$$

After that, we compute the  $\delta$  terms for the earlier layers:

$$\delta^{(3)} = (\Theta^{(3)})^T \cdot \delta^{(4)} \circ g'(z^{(3)}).$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \circ g'(z^{(2)}).$$

$\circ$  element-wise multiplication.

The derivative  $g'(z^{(3)})$  is basically a vector of one:

$$g'(z^{(3)}) = \underbrace{a^{(3)}}_{\text{vector}} * \underbrace{(I - a^{(3)})}_{\text{vector}}$$

$$g'(z^{(2)}) = a^{(2)} * (1 - a^{(2)}).$$

→ there is no  $f^{(1)}$  since they are the features we observe in our training set

→ The name back propagation comes from the fact we start in the output layer,

calculate the error  $f^{(4)}$  ~~at~~ there and go back to the 3rd layer & where we compute  $\delta^{(3)}$  and go back to layer 2 and calculate  $\delta^{(2)}$  in such a way it is propagating the error from layer  $(l+1)$  to  $(l)$  and from  $(l)$  to  $(l-1)$  and so on.

Finally 
$$\boxed{\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)}} \text{ when ignoring the regularization term } (\lambda=0).$$

Now, for a larger training set  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ , set

$$\Delta_{ij}^{(l)} = 0 \quad \forall i, j, l.$$

→ will be used as accumulators for computing  $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$ .

For  $i=1$  to  $m\}$  // on each iteration, we work with one example.

$$\left| \text{Set } a^{(1)} = x^{(i)} \right.$$

Perform forward propagation to compute  $a^{(l)}$  for  $l=2, 3, \dots, L$ .

Using  $y^{(i)}$ , compute  $\delta^{(L)} = a^{(L)} - y^{(i)}$  // computing error for this example for the output layer.

Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$  used to compute the partial derivatives in the end  $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$ .

$$\boxed{\Delta_{ij}^{(l)} \leftarrow \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}}$$

$$\left\{ \begin{array}{l} \Delta_{ij}^{(l)} \leftarrow \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)} \text{ if } j \neq 0 \\ \Delta_{0j}^{(l)} \leftarrow \frac{1}{m} \cdot \Delta_{0j}^{(l)} \text{ if } j=0 \end{array} \right. \text{ (bias term)} \quad \text{if } j \neq 0$$

if we vectorize,

$$\boxed{\Delta^{(l)} \leftarrow \Delta^{(l)} + \delta^{(l+1)} \cdot (a^{(l)})^T}$$

is exactly

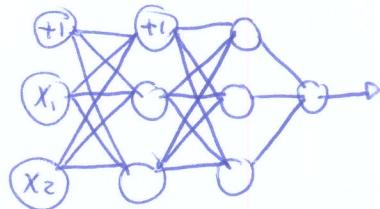
$$\boxed{\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = \Delta_{ij}^{(l)}}$$

practice

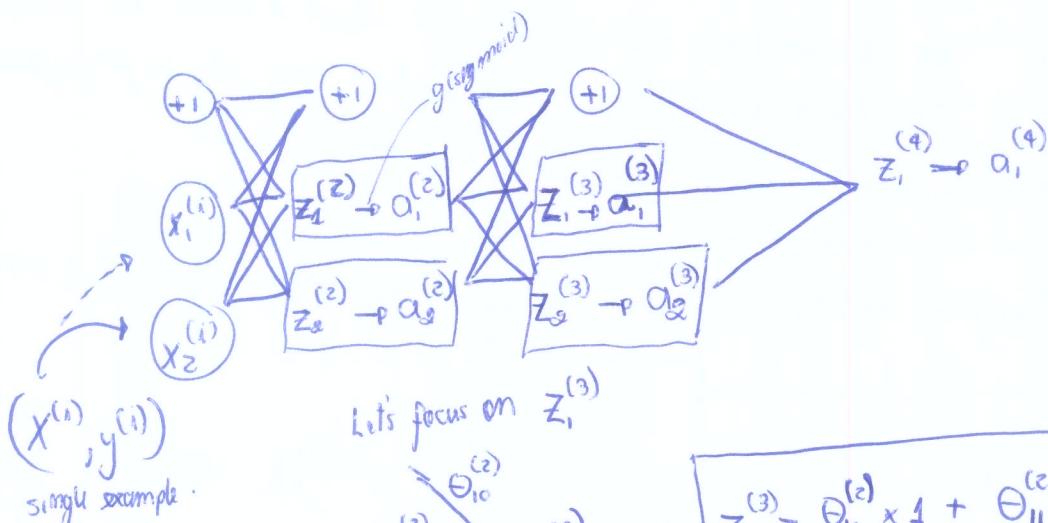
## Backpropagation Intuition

Mechanical steps.

Let's see again forward propagation:



units & sum sum 1 unit  
not counting bias



$$\text{So } z_1^{(3)} = \Theta_{10}^{(2)} \cdot 1 + \Theta_{11}^{(2)} \cdot a_1^{(2)} + \Theta_{12}^{(2)} \cdot a_2^{(2)}$$

this is forward propagation.

What does Backpropagation do?

Recall its formulation:

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log(f_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-f_\theta(x^{(i)})) \right] + \underbrace{\frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{j=1}^{S_l} \sum_{k=1}^{S_{l+1}} (\Theta_{jk}^{(l)})^2}_{\text{Regularization.}}$$

focusing on a single example  $(x^{(i)}, y^{(i)})$ , the case of ① output unit and  $(\lambda=0)$

$$\text{Cost}(i) = y^{(i)} \log f_\theta(x^{(i)}) + (1-y^{(i)}) \log \cancel{(1-f_\theta(x^{(i)}))}$$

$$\text{think of } \text{Cost}(i) \approx f_\theta(x^{(i)}) - y^{(i)}_i$$

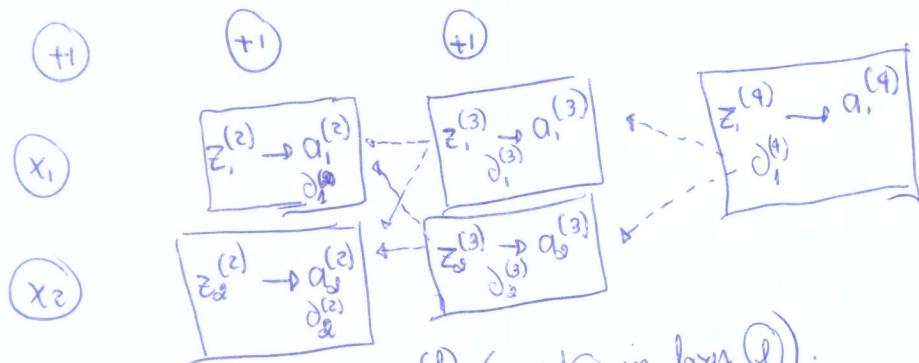
How is the NN on example  $i$ ?

the correct form is

$$\log(1-f_\theta(x^{(i)}))$$

## Forward Propagation

15



$f_j^{(l)}$  = error of cost for  $a_j^{(l)}$  (unit  $j$  in layer  $l$ )

Formally,  $\delta_j^{(l)} = \frac{\partial \text{cost}(i)}{\partial z_j^{(l)}}$  change how affects  $\text{cost}(i)$  (for  $j \geq 0$ ) where

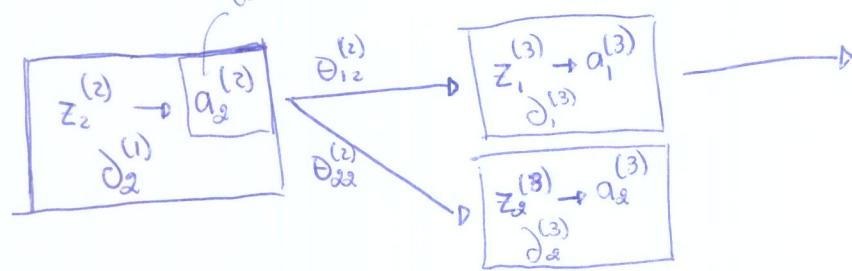
$$\text{cost}(i) = y^{(i)} \cdot \log f_0(x^{(i)}) + (1 - y^{(i)}) \cdot \log (1 - f_0(x^{(i)}))$$

\* here and \* here

Let's see it:

$$\text{for } \delta_1^{(4)} = y^{(1)} - a_1^{(4)} \quad \text{output}$$

how to get  $\delta_2^{(2)} = ?$



$$\delta_2^{(2)} = \Theta_{12}^{(2)} \delta_1^{(3)} + \Theta_{22}^{(2)} \cdot \delta_2^{(3)}$$

and

$$\delta_2^{(3)} = \Theta_{12}^{(3)} \cdot \delta_1^{(4)}$$

— x —

15  
math

## Implementation Note

Example  $s_1 = 10$  units

$$s_2 = 10$$

$s_3 = 1$  output unit

$$\Theta^{(1)} \in \mathbb{R}^{10 \times 11}$$

$$D^{(1)} \in \mathbb{R}^{10 \times 11} \text{ derivatives}$$

$$\Theta^{(2)} \in \mathbb{R}^{10 \times 11}$$

$$D^{(2)} \in \mathbb{R}^{10 \times 11}$$

$$\Theta^{(3)} \in \mathbb{R}^{1 \times 11}$$

$$D^{(3)} \in \mathbb{R}^{1 \times 11}$$

If we want to use this with more advanced ~~#~~ techniques of optimization,  
we may need to deal with them as ~~matrices~~ instead of as ~~matrices~~ vectors.

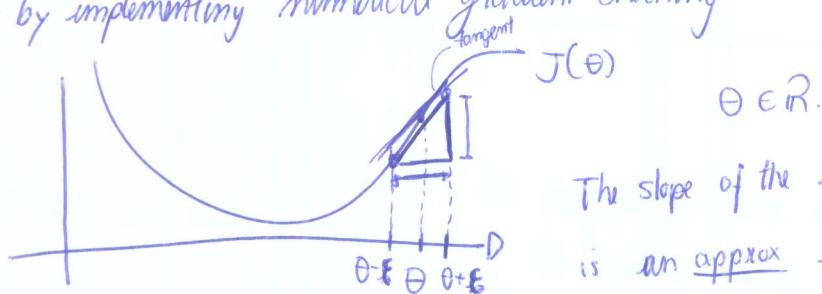
E.g., when optimizing in Octave.



## Gradient Checking

How can we check the Backpropagation algorithm along with its optimization function for finding the parameters (e.g., gradient descent) are actually working?

~ We start by implementing numerical gradient checking.



$$\theta \in \mathbb{R}$$

The slope of the line  $(J(\theta + \epsilon), (\theta + \epsilon); J(\theta - \epsilon), \theta - \epsilon)$  is an approx for the derivative of  $J(t)$ .

$$\text{slope} = \frac{\text{vertical height}}{\text{horizontal width}} = \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}.$$

Therefore  $\frac{\partial}{\partial \theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$  } two sided difference.

$\epsilon = 10^{-4}$  (the smaller the  $\epsilon$ , the closer to the actual derivative).

In a more general case,  $\theta \in \mathbb{R}^n$  (e.g.,  $\theta$  is "unrolled" version of  $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$ )  
unrolled matrix  $d \times d$   $\begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$ .

$$\theta = [\theta_1, \theta_2, \dots, \theta_n]^T.$$

we can do:

$$\frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\underline{\theta_1 + \epsilon}, \theta_2, \theta_3, \dots, \theta_n) - J(\underline{\theta_1 - \epsilon}, \theta_2, \dots, \theta_n)}{2\epsilon}.$$

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \underline{\theta_2 + \epsilon}, \theta_3, \dots, \theta_n) - J(\theta_1, \underline{\theta_2 - \epsilon}, \theta_3, \dots, \theta_n)}{2\epsilon}.$$

$$\frac{\partial}{\partial \theta_m} J(\theta) \approx \frac{J(\theta_1, \theta_2, \dots, \underline{\theta_m + \epsilon}) - J(\theta_1, \theta_2, \dots, \underline{\theta_m - \epsilon})}{2\epsilon}.$$

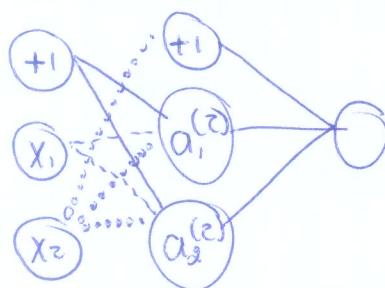
These equations allow us to numerically estimate the derivatives with respect to any parameter  $\theta_i$  and then can be used to double check the derivatives calculated by back propagation.

~ Backpropagation is very efficient compared to gradient numerical checking.  
 Therefore we should use it just in one or a few iterations to ~~do~~ double check backpropagation results.

### Random Initialization

We need to pick some initial value for  $\Theta$ .  
 Initializing  $\Theta$  with zeros worked OK for logistic regression but is ~~doen't~~ for NNs.  
 The reason is that the network will be redundant.

$a_1^{(2)} = a_2^{(2)}$  and  $\delta_1^{(2)} = \delta_2^{(2)}$ .  
 After each update, parameters corresponding to inputs going into each of two hidden units are identical



$a_1^{(2)} = a_2^{(2)}$  if  $\Theta = 0$  for all elements in the beginning.

This is known as the problem of symmetric weights. For breaking the symmetry, we must initialize  $\Theta_{ij}^{(l)}$  to a random value in  $[-\epsilon, \epsilon]$  ( $\epsilon$  means  $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$ ) This is not  $\Theta$  of before.

### Wrapping up

The first thing we need to do when training a NN, is to select a NN architecture (connectivity pattern). 

How to choose this?

- ① input units: dimensionality of the problem  $X^{(1)}$ .
- ② output units: number of classes.

③ for the internal (hidden layers), a reasonable default is one hidden layer. But if we are using more than one hidden layer, it is interesting to have the same number of units in all layers (the more, the better, normally). hidden layers

Note : } the more hidden layers, the better at the cost of much more computation.  
} nbr of units in hidden layer normally is comparable to that in the input layer.

## Training a neural network

① Initialize weights randomly.

② Implement forward propagation to obtain  $f_{\theta}(x^{(i)}) \approx x^{(i)}$ .

③ Compute cost function  $J(\theta)$ .

④ Backpropagation to compute partial derivatives  $\frac{\partial}{\partial \theta_j^{(l)}} J(\theta)$ .

Example for  $i=1$  to  $(m)$  } <sup>more examples</sup> //  $x^{(1)}, y^{(1)}; \dots; x^{(m)}, y^{(m)}$ .

| Perform forward prop. and backprop using  $x^{(i)}, y^{(i)}$ .

| Get activation's  $a^{(1)}$  and delta terms  $\delta^{(l)}$  for  $l=2, \dots, L$ .

|  $\Delta^{(1)} \leftarrow \Delta^{(1)} + \delta^{(L+1)} (a^{(1)})^T$ . // accumulators having, in the end, the partial derivative  
| }  $\frac{\partial}{\partial \theta_j^{(1)}} J(\theta)$

Compute  $\frac{\partial}{\partial \theta_j^{(1)}} J(\theta)$

⑤ Use gradient checking to compare  $\frac{\partial}{\partial \theta_j^{(1)}} J(\theta)$  computed in ④ with a numerical estimate of gradient of  $J(\theta)$ .

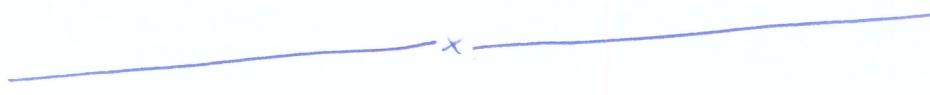
⑥ Disable gradient checking.

⑦ Use Gradient Descent or other opt method with backprop to minimize  $J(\theta)$  as a function of params  $\theta$ .

| As we have  $J(\theta)$  as a non-convex function, we may ~~not~~ converge to local minima.

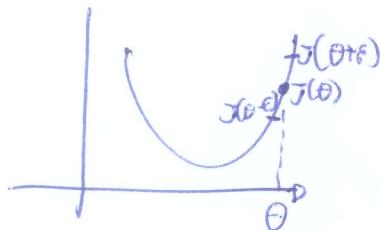
Answers

} Backprop  $\rightarrow$  calculates the derivatives (direction of a step).  
Gradient Descent  $\rightarrow$  take little baby steps down <sup>the gradient</sup> hill.



### Last notes on Gradient checking

- Always check the first results of Gradient Descent using gradient checking.



- Gradient checking is too slow. Do not use it for learning, just to check if it works in the first iteration.

