

Representation

LR might not be the suitable solution for large problems with non-linear hypothesis

Why NN? not enough with LR?

the problem becomes too large \therefore

when fitting or classifying non-linear problems the number of features might increase a lot:

$$x_1, \dots, x_{100}$$

original
features
(n=100)

\rightarrow quadratic features $\rightarrow x_1^2, x_1 x_2, x_1 x_3, \dots, x_1 x_{100}, x_2^2, \dots$
 $O(n^2) \approx 5000$ features

\rightarrow cubic features $\rightarrow x_1^3, x_1 x_2^2, x_1^2 x_2, \dots$
 $\approx 170,000$ features
 $O(n^3)$

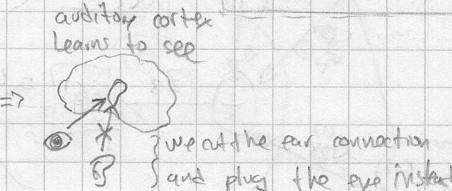
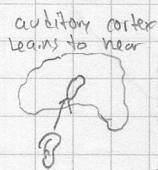


Images is another example where LR is not enough, because each pixel is treated as a feature...

$$\begin{array}{l} [50 \times 50] \rightarrow 2500 \text{ features} \\ \text{(7500 if RGB)} \end{array} \rightarrow n = 2500 \rightarrow \text{quadratic features} \approx 3 \text{ million } \therefore$$

Neurons and brain

The brain is based on the "one learning algorithm" hypothesis because part is prepared to learn from any input signals.

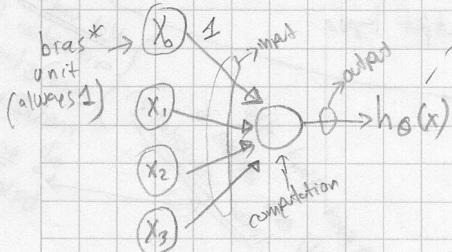


? but, LR also follows this principle, right? so this "one learning algorithm" description could also apply to LR...

this is the idea of NN,
 → one algorithm to fit all problems.

Model Representation

Single neuron model



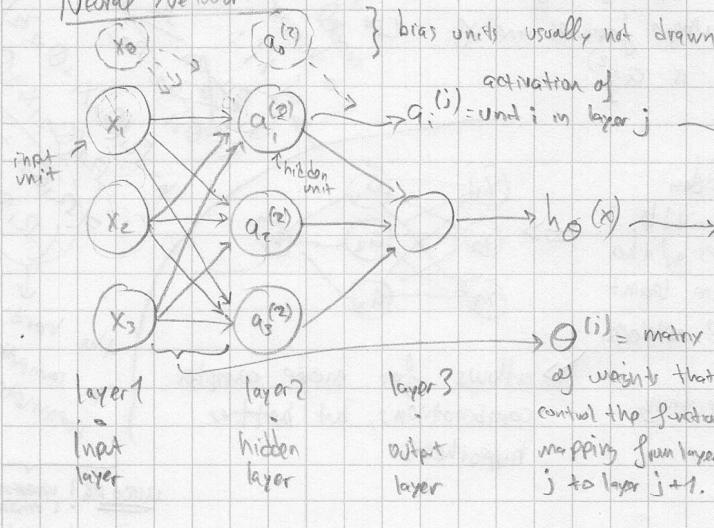
$$h_0(x) = \frac{1}{1 + e^{-\Theta^T x}} \quad \text{where } x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \text{ and } \Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

* the bias unit is sometimes not drawn ($x_0 = 1 \dots$)

I guess it is a sigmoid function for model activation of that unit.

? Is Logistic Regression a single neuron?

Neural Networks



$$\begin{aligned} a_1^{(2)} &= g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3) \\ a_2^{(2)} &= g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3) \\ &\vdots \\ a_3^{(2)} &= g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3) \end{aligned}$$

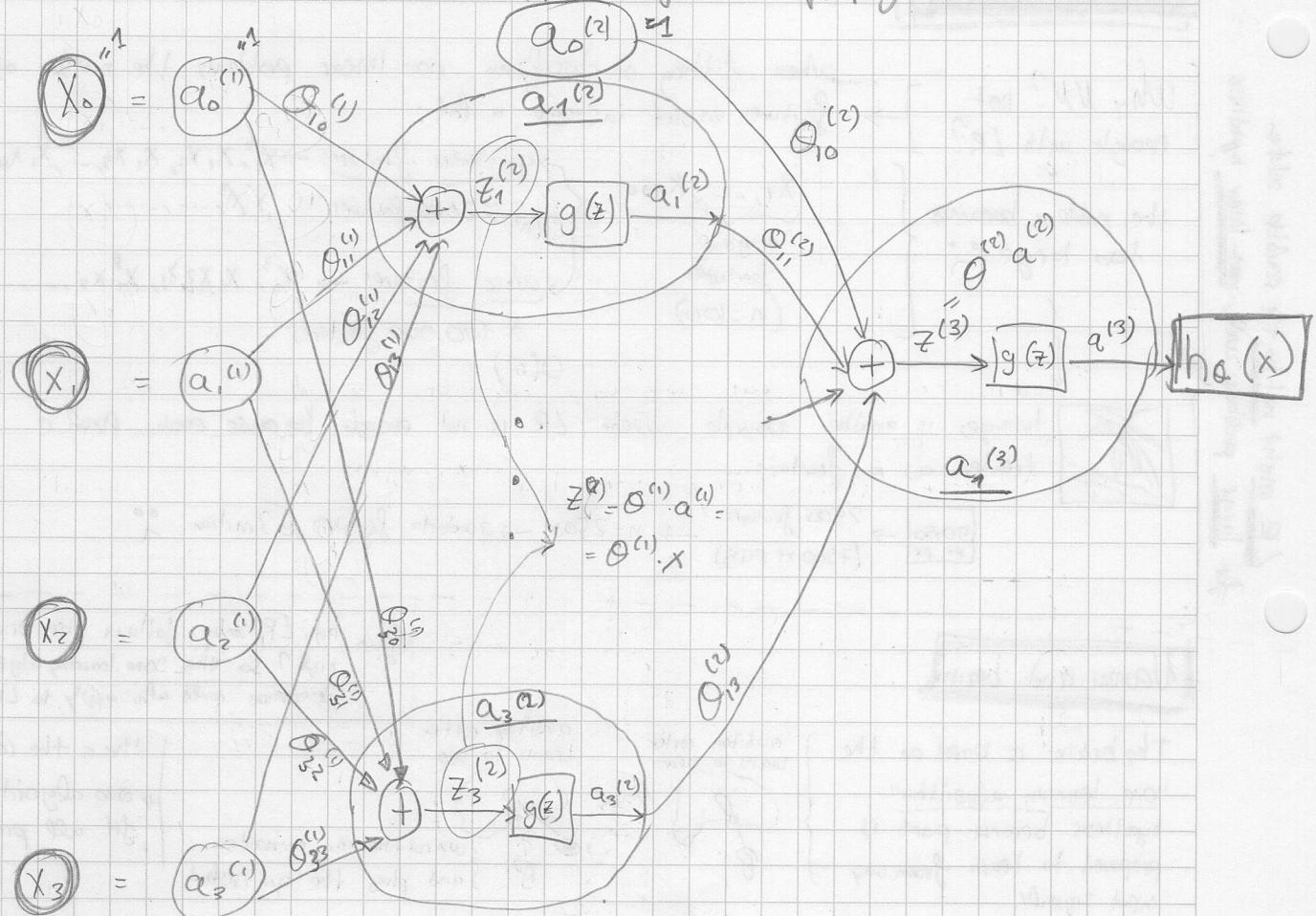
$\Theta^{(i)}$ = matrix of weights that control the function mapping from layer j to layer $j+1$.

dimension of $\Theta^{(i)}$ is $(\dim g(\Theta^{(1)}), 3 \times (3+1))$
 \downarrow
 $S_{j+1} \times (S_j + 1)$
 $\# \text{units in layer } j+1$
 \downarrow
 $\# \text{units in layer } j$
 $\not= \text{counting } a_0$
 $\not= \text{counting } x_0$

A more detailed diagram

→ to show vectorization implementation

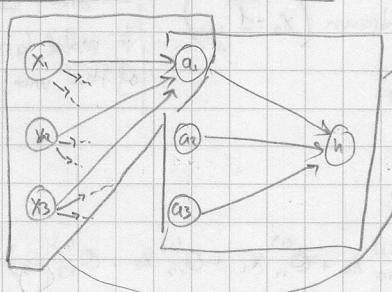
→ to explain 'forward propagation'



Forward Propagation → start with activation of input unit and propagate to hidden layer (computing activations of hidden items), and forward propagate to compute the activation of the output layer

Interpretation of steps

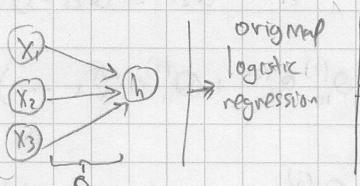
Using $g(z)$ helps with the idea of activating each layer element



this is just applying logistic regression where features are $a_1^{(2)}$

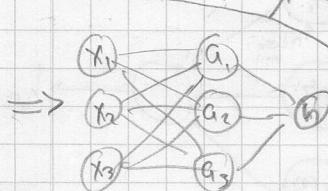
this is just applying logistic regression where features are x_i and output is $a_1^{(2)}$

I interpret NN that we can obtain a hierarchical linear combination of non-linear behaviour!



original logistic regression \Rightarrow we add a hidden layer in the middle and the features of this hidden layer are learnt by the neural network

this is the nice thing about NN



\Rightarrow allows for more complex combinations, but better hypothesis

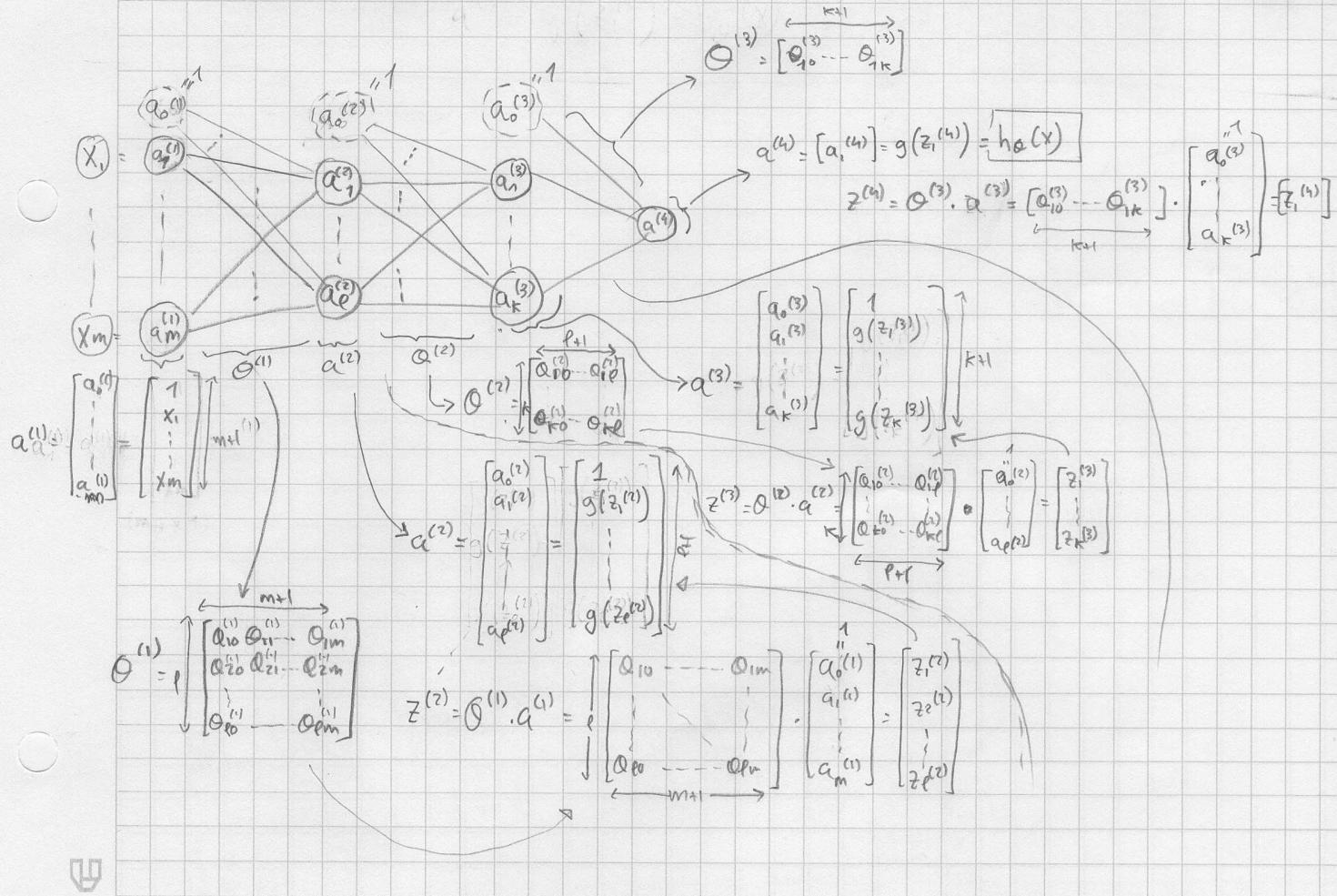
$$h = P_1(O_1 x_1 + O_2 x_2) + P_2(O_3 x_1 + O_4 x_2) + P_3(O_5 x_1 + O_6 x_2) + P_4(O_7 x_1 + O_8 x_2)$$

$$= P_1(O_1 x_1 + O_2 x_2) + P_2(O_3 x_1 + O_4 x_2) + P_3(O_5 x_1 + O_6 x_2) + P_4(O_7 x_1 + O_8 x_2)$$

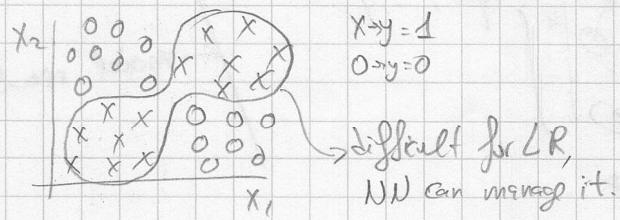
$$= P_1(O_1 x_1 + O_2 x_2) + P_2(O_3 x_1 + O_4 x_2) + P_3(O_5 x_1 + O_6 x_2) + P_4(O_7 x_1 + O_8 x_2)$$

the hard part is to compute the right parameter values.

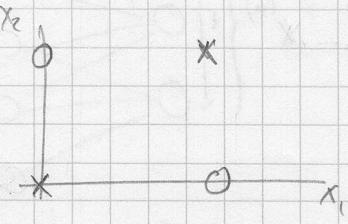
we'll see in the next slide to consider the activation functions at each $a_j^{(i)}$



Example and intuition



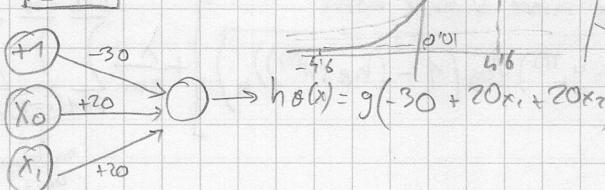
simplification
 \rightsquigarrow



$X_1 \text{ XNOR } X_2$
(1 only if $X_1 = X_2$)

He reviews how to model different operators using NN, but in all cases he gives the parameters, so it is more an exercise to understand how NN compute values (forward pass) rather than understanding how to come up with Θ .

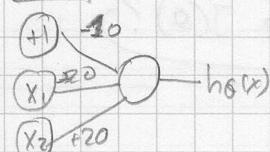
[AND]



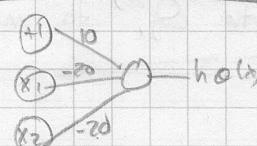
x_1	x_2	$h_\theta(x)$
0	0	$g(-30) \approx 0$
0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
1	1	$g(10) \approx 1$

remember this is $P(y=1|x_1, x_2)$

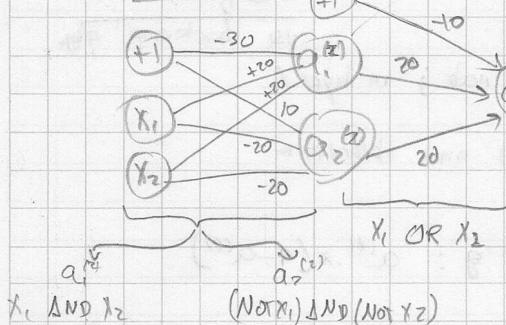
[OR]



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



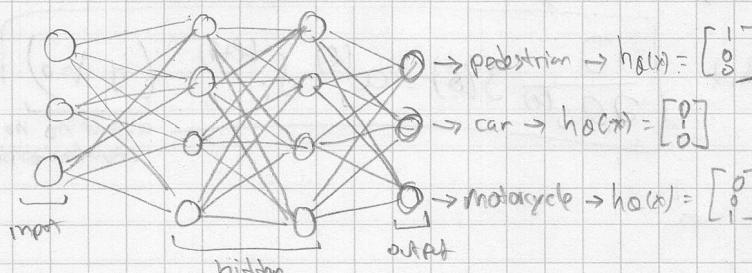
$X_1 \text{ XNOR } X_2$



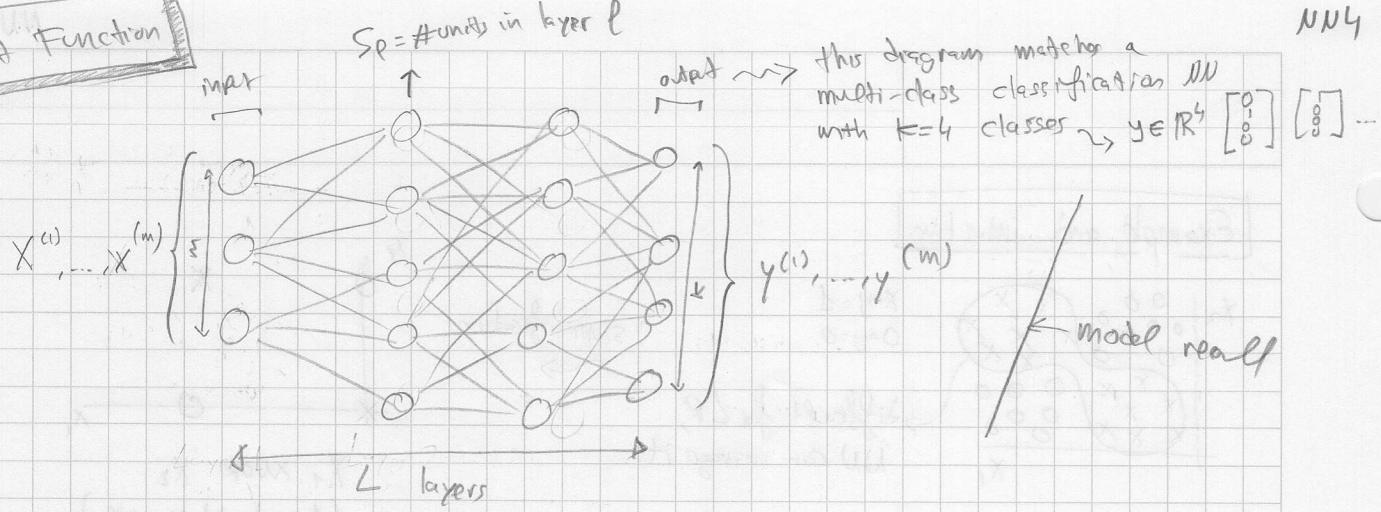
x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_\theta(x)$
0	0	0	0	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

Multiclass classification

(same idea as in one-vs-all)



Cost Function



The cost function will be a generalization of the log. Reg. cost function:

$$\text{Logistic Regression} \leadsto J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1-h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \left[\sum_{k=1}^K \left[\begin{array}{c} \text{cost of single sample} \\ \text{2 single output} \end{array} \right] \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{j=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2 \right]$$

Annotations for the equation:

- all m samples
- all K outputs
- i th sample
- k th output observed
- k th output predicted
- from layer l to layer $l+1$
- no bias unit
- layer $l+1$

but how to compute $\theta_{ji}^{(l)}$ to minimize $J(\theta)$?

Back propagation algorithm

To minimize $J(\theta)$, we need $J(\theta)$ and the partial derivatives already have it!

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

use backprop. alg.

The algorithm is based on computing $\delta_j^{(l)}$ \Rightarrow error of node j in layer l .

- $\delta_j^{(l)} = a_j^{(l)} - y_j \quad \Rightarrow$ difference between hypothesis and observation
- $\delta^{(l)} = a^{(l)} - y \quad \leftarrow$ vectorized
- $\delta^{(l)} = (\theta^{(l)})^T \delta^{(l+1)}$ $\underbrace{g'(z^{(l)})}_{\substack{\text{elements} \\ \text{wise}}} \quad \Rightarrow$ derivative of g : $a^{(l)} \cdot (1-a^{(l)})$

(no $\delta^{(1)}$ \rightarrow layer 1 is observations, no error...)

Having $\delta_j^{(l)}$, it can be proved that

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)} \quad (if \lambda=0)$$

assuming no regularization

actually, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$,
so it measures how the cost changes when we change $z_j^{(l)}$, and the change z by changing θ ...

The algorithm

Set $\Delta_{ij}^{(l)} = 0$

loop all training set

For $i=1$ to m ← introduce input observation

Set $a^{(1)} = x^{(i)}$

Forward propagation to compute $a^{(2)}, a^{(3)}, \dots, a^{(L)}$ ← but with which $Q_{ij}^{(k)}$???

Compute $f^{(L)} = a^{(L)} - y^{(i)}$

Backpropagation to compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}, \delta^{(1)}$

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ ← accumulate the partial derivative terms in δ

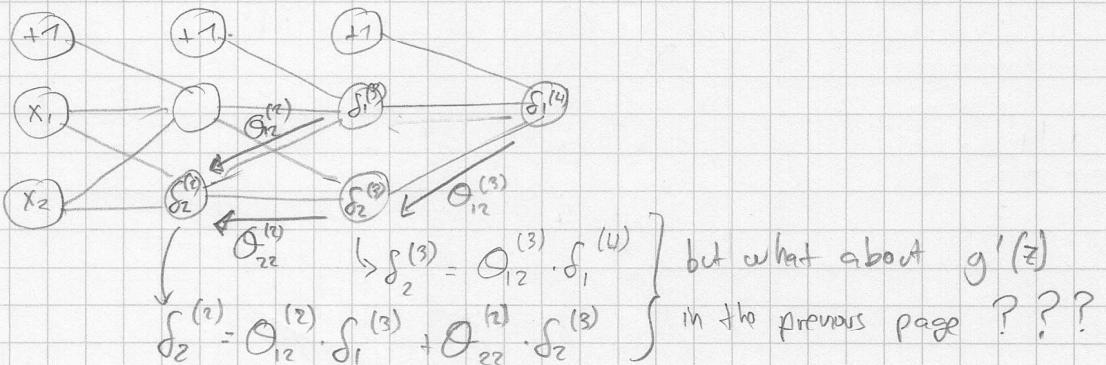
$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda Q_{ij}^{(l)}, \text{ if } j \neq 0$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}, \text{ if } j = 0$$

$$\rightarrow D_{ij}^{(l)} = \frac{\partial}{\partial Q_{ij}^{(l)}} J(Q) \rightarrow \text{so we have the partial derivative of the cost function !!!}$$

The intuition

It is just like forward propagation, but backwards...



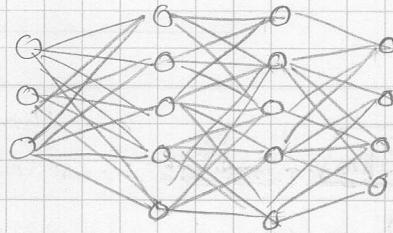
A. Pick a network architecture

input units = dimensions of X

output units = number of classes

hidden layers = Default 1 (if > 1 , have same number of units per layer)

hidden units = slightly larger than # input units (2-3-4 times is ok)



B. Train the neural network

1. Randomly initialize weights

2. Forward Propagation to get $h_{\theta}(x^{(i)})$ for any $x^{(i)}$

3. Implement code to compute $J(\theta)$

4. Back prop to compute $\frac{\partial}{\partial \theta_j^{(k)}} J(\theta)$

(5. Use gradient checking to compare $\frac{\partial}{\partial \theta_j}$ from back prop vs. numerical estimate of gradient of $J(\theta)$)

6. Use gradient descent or advanced optimization to minimize $J(\theta)$ as a function of θ .