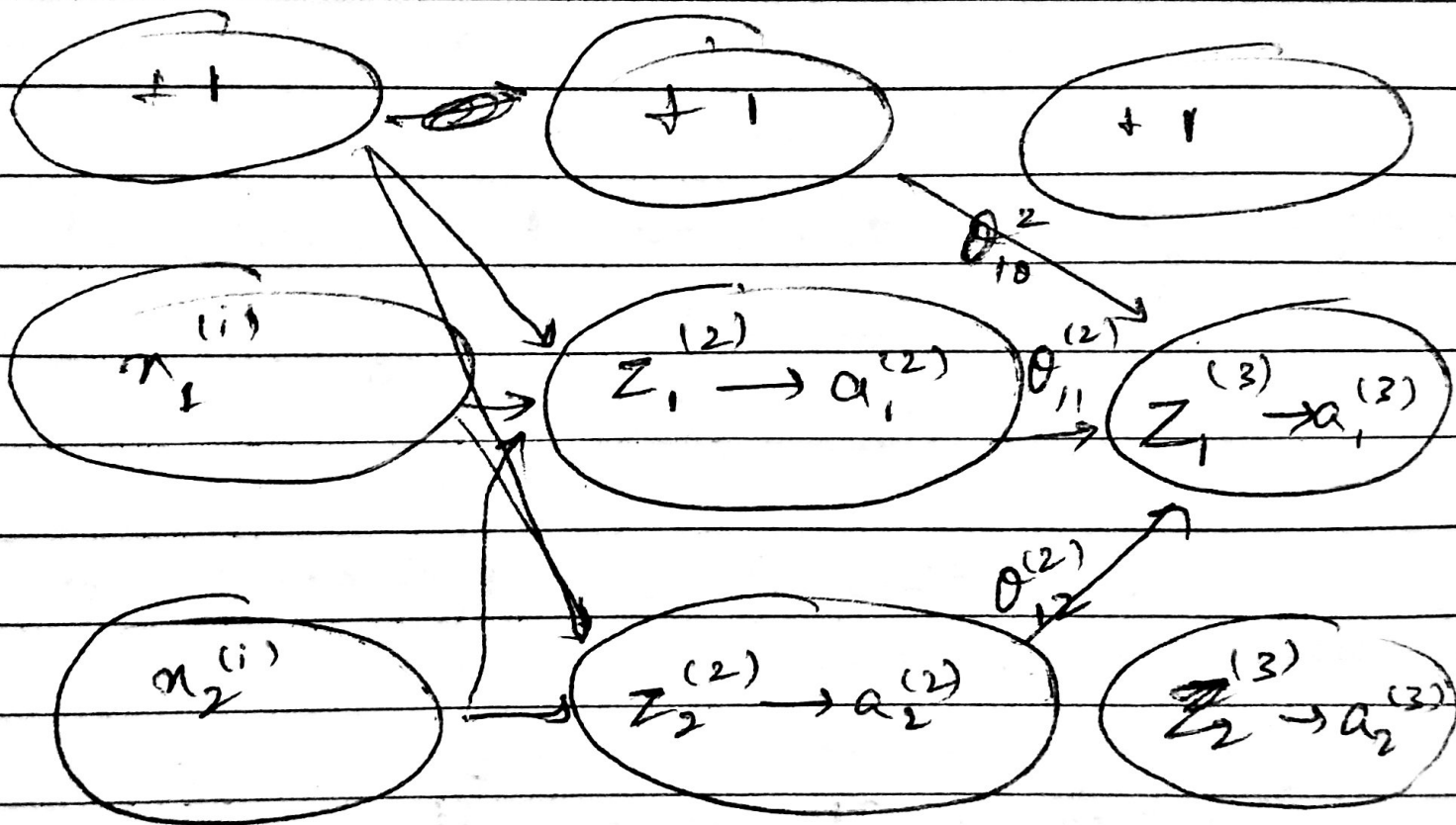


Forward propagation



\therefore for
$$z_1^{(3)} = \theta_{10}^{(3)} \cdot 1 + \theta_{11}^{(3)} a_1^{(2)} + \theta_{12}^{(3)} a_2^{(2)}$$

Unrolling Parameters

Normally, we have

$$\text{function } [J, \text{Val}, \text{gradients}] = \text{costfunction}(\text{theta})$$

\uparrow cost func \uparrow parameters (vector) \uparrow vector

optTheta = fminunc (@ costfunction, initialTheta, options)

\uparrow vector

however, in neural network,

$\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ } all matrices
 $D^{(1)}, D^{(2)}, D^{(3)}$

suppose, $S_1 = 10, S_2 = 10, S_3 = 1$

$$\Theta^{(1)} \in \mathbb{R}^{10 \times 11} \quad \Theta^{(2)} \in \mathbb{R}^{10 \times 11} \quad \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$$

unroll: theta_vec = [Theta1(:); Theta2(:); Theta3(:)]
D_vec = [D1(:); D2(:); D3(:)];

for getting back,
theta1 = reshape(theta_vec(1:110), 10, 11)
theta2 = reshape(theta_vec(111:220), 10, 11)
theta3 = reshape(theta_vec(221:231), 1, 11)

Now,
for function,

function[jval, gradientVec] = costFunction(thetaVec)

↑
vector

→ From thetaVec,
get $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$

→ use forward prop/back prop to
compute $D^{(1)}, D^{(2)}, D^{(3)}$ & $J(\theta)$

→ unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradients

GRADIENT CHECKING (Bug fixer)

$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

$$\boxed{\epsilon \approx 10^{-4}}$$

for implementation,

$$\text{gradApprox} = \frac{J(\text{theta} + \epsilon \text{vector}) - J(\text{theta} - \epsilon \text{vector})}{2 * \epsilon \text{vector}}$$

Def 2ⁿ

$$\theta = \theta_1, \theta_2, \theta_3, \dots, \theta_n$$

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

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

Implementation

for $i = 1:n$

thetaPlus = theta;

thetaPlus(i) = thetaPlus(i) + epsilon;

thetaMinus = theta;

thetaMinus(i) = thetaMinus(i) - epsilon;

$$gradApprox(i) = (J(\text{thetaPlus}) - J(\text{thetaMinus})) / (2 * \epsilon)$$

end;

check gradApprox \approx Dvec

↑ from
backpropagation

→ Turn off gradient checking after you have checked the Dvec & gradApprox because gradient checking is very slow.

RANDOM INITIALIZATION

Some random θ

$optimizer = fminunc(@costFunction, initialTheta, options)$

we usually set,

$initialTheta = zeros(n, 1)$

for regression, ~~for~~ it should not be used for Neural Networks.

for neural network, we randomly initialize weights in between $-G$ & G
i.e., $-G \leq \theta_{ij}^{(1)} \leq G$

$Theta1 = rand(10, 11) * (2 * INIT_EPSILON) - INIT_EPSILON$

$Theta2 = rand(4, 11) * (2 * INIT_EPSILON) - INIT_EPSILON$

How to

- 1) We start by picking a network architecture (connectivity pattern b/w neurons)

i/p is $x(i)$

o/p is no. of classes. if y takes values as $y \in \{1, 2, 3 \dots 10\}$

then, you have a vector of dimension ten.

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Reasonable default is take 1 hidden layer. or if you take more, same no. of units in each layer.