

\Rightarrow Backpropagation ~~intuition~~ intuition

\rightarrow Backpropagation is actually doing something very similar to forward propagation, but backwards!

\rightarrow The cost funcⁿ for our NN is:

$$J(\theta) = -\frac{1}{m} \sum_{t=1}^m \sum_{k=1}^K [y_k^t \log(h_\theta(x^t))_k + (1-y_k^t) \log(1-h_\theta(x^t))_k] \\ + \frac{\lambda}{2m} \sum_{L=1}^{L-1} \sum_{i=1}^{d_L} \sum_{j=1}^{d_{L+1}} (\theta_{ji})^2$$

\rightarrow If we consider simple non-muticlass classification ($K=1$),

& disregard regularization, the cost is computed with:

Cost(x)

$$J(\theta) = y \log(h_\theta(x)) + (1-y) \log(1-h_\theta(x))$$

for one training example.

* We can think about a δ term on a unit as the "error" of cost for the activation value associated with a unit.

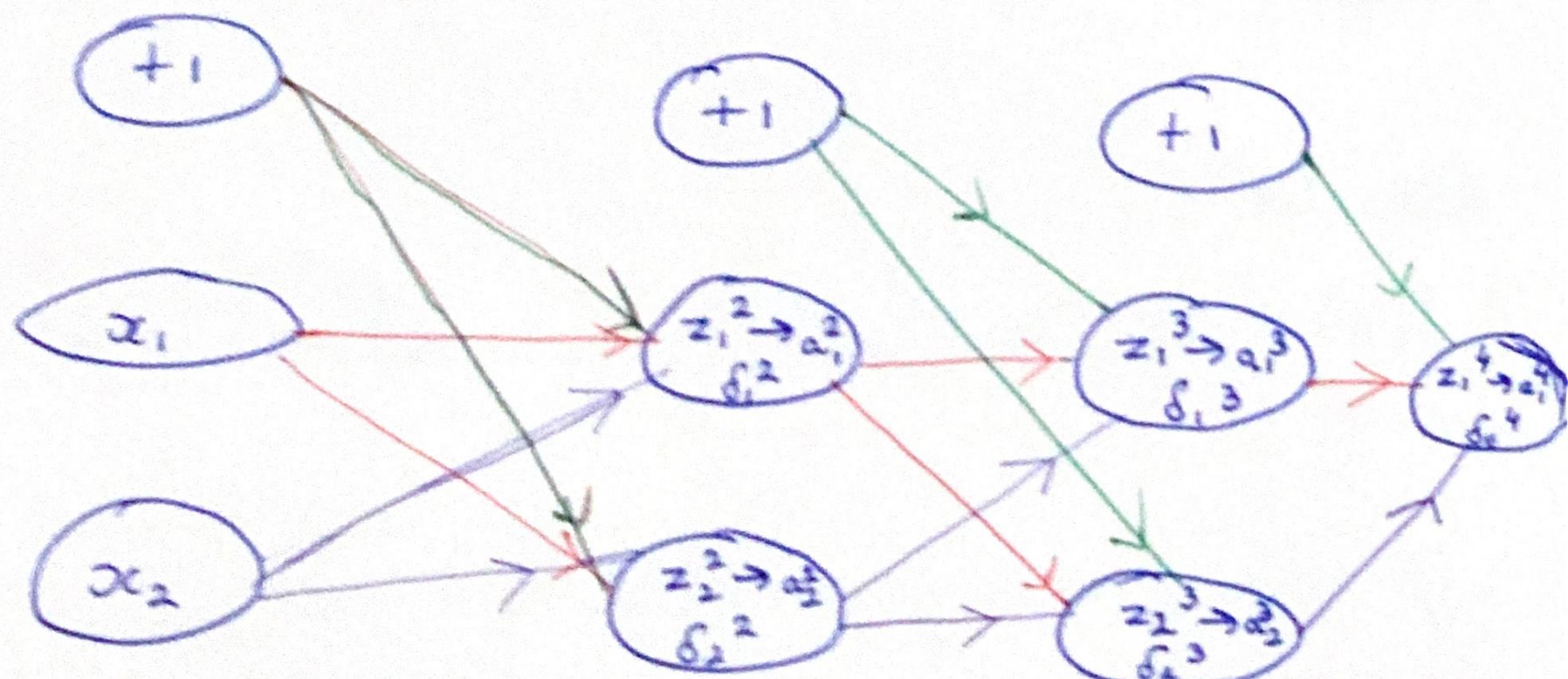
→ More formally,

$$\delta_j^L = \frac{\partial \text{Cost}(\theta, \alpha)}{\partial z_j^L}$$

* For the output layer $\boxed{\delta^L = a - y}$

We propagate these values backwards.

e.g.:



$$\delta^L = y - a^L$$

$$\therefore \delta_1^2 = \Theta_{11}^2 \delta^L$$

$$\delta_2^3 = \theta_{12}^3 \delta_1^4$$

→ Continuing to prev. layers:

$$\delta_1^2 = \theta_{11}^2 \delta_1^3 + \theta_{21}^2 \delta_2^3$$

$$\delta_2^2 = \theta_{12}^2 \delta_1^3 + \theta_{22}^2 \delta_2^3$$

→ So, backpropagation order:

* Calculates the δ ,

* And these δ values are the weighted sum of the next layer's delta values.

* δ corresponds to nodes Δ corresponds to weights

⇒ Implementation note: Unrolling/Flattening parameters


→ With neural networks, we are working with sets of matrices :

$$\Theta^1, \Theta^2, \Theta^3, \dots$$

$$D^1, D^2, D^3, \dots$$

→ In order to use optimizing func's such as fminunc,
we'll want to "unroll" or "flatten" all the elements into a one long vector :

$$\text{thetaVector} = [\Theta^1(:); \Theta^2(:); \Theta^3(:); \dots]$$

$$dVector = [D^1(:); D^2(:); D^3(:); \dots]$$

→ We can get back off our original matrices by using a bit of slicing & "reshape" func.

e.g.: if

$\dim(\Theta_{\text{eta}1}) = 10 \times 11$
 $\dim(\Theta_{\text{eta}2}) = 10 \times 11$
 $\dim(\Theta_{\text{eta}3}) = 1 \times 11$

then we can get back our original matrices as:

$\Theta_{\text{eta}1} = \text{reshape}(-\theta_{\text{etaVector}}(1:110), 10, 11)$

$\Theta_{\text{eta}2} = \text{reshape}(-\theta_{\text{etaVector}}(111:220), 10, 11)$

$\Theta_{\text{eta}3} = \text{reshape}(-\theta_{\text{etaVector}}(221:231), 1, 11)$

→ By doing this, our learning algs for the above eg. becomes:

Having initial parameters $\theta^1, \theta^2, \theta^3$

Unroll / Flatten to get "initial Theta" to pass to
`fminunc (@CostFunction, initialTheta, options)`

function [fval, gradientVec] = costFunction(thetaVec)

From thetaVec, get $\theta^1, \theta^2, \theta^3$ back using reshape

Use forward prop & backward prop to compute
 $\delta^1, \delta^2, \delta^3$ & $J(\theta)$

Unroll / Flatten $\delta^1, \delta^2, \delta^3$ to get gradient Vec

⇒ Gradient checking



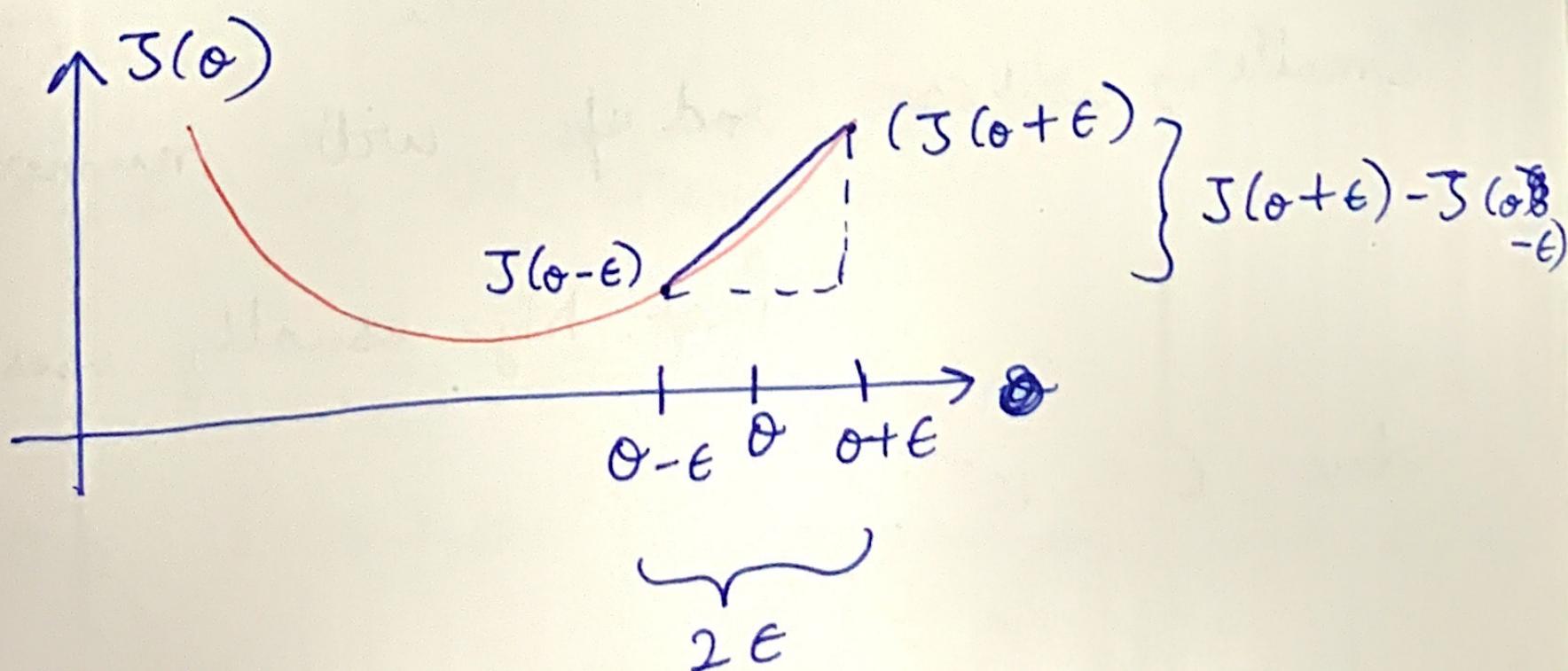
→ Backpropagation can be messy! It is easy to miss backprop some imp. details and while implementing it is easier to face bugs.

e.g.:

This might mean that

→ Gradient checking will assure that our backpropagation works as intended.

→ Let $J(\theta)$ be a single valued function only dependent on θ . Supp. $J(\theta)$ looks like:



If $\epsilon \rightarrow 0$ $\therefore \boxed{\frac{d J(\theta)}{d \theta} \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}}$

* This is the 2 sided difference δ is more accurate than a 1 sided difference $\left[\frac{\partial J(\theta)}{\partial \theta} = \frac{J(\theta + \epsilon) - J(\theta)}{\epsilon} \right]$

* With multiple theta matrices, we can approximate the derivative wrt θ_j as follows:

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

* A good small value for ϵ guarantees the method above to become true. If the value be much smaller, will end up with numerical problems.

bsf. Ng usually uses the value $\epsilon = 10^{-4}$.

* Using the prev. algs, we can manually (or rather laboriously) calculate the value of $\frac{\partial J(\theta)}{\partial \theta_{ij}}$ for each of every weight in ~~each~~ our network!

We can see why we need backpropagation's ~~as backpropagation~~ why it is so efficient.

As backpropagation alg
Backpropagation works by computing the gradient of the loss funcⁿ wrt each weight by the chain rule.

By iterating backwards from the last layer, it computes the gradient one layer at a time using the gradients of the next layer (Rev. wrt iterating backwards!), it removes the ~~redundant~~ redundant calculations of the intermediate terms.

- * Backpropagation can be a clear example of dynamic programming!
- * Manually calculating gradient can be slow but it helps us in stress testing our backpropagation algs , to check whether backprop is calculating the right gradient values or not.
- * We can implement this "slow" computation as follows :

Let there be m no. of parameters in our cost function i.e. $J(\theta) = J(\theta_1, \dots, \theta_m)$

$$\epsilon \text{psilon} = 1e^{-4};$$

for $i = 1 : m,$

$\theta_{\text{plus}} = \theta;$

$\theta_{\text{plus}}(i) + = \epsilon \text{psilon};$

$\theta_{\text{minus}} = \theta;$

$\theta_{\text{minus}}(i) - = \epsilon \text{psilon};$

$\text{grad } \theta_{\text{approx}}(i) = \left[\frac{J(\theta_{\text{plus}}) - J(\theta_{\text{minus}})}{2 * \epsilon \text{psilon}} \right];$

end

→ We earlier saw how we can compute the "D" matrix using backpropagation.

So, once we compute our gradApprox vector, we can check that whether $D = \text{grad Approx}$ or not.

→ Once, you've verified that backprop is working correctly; don't forget to turn off the code for computing gradApprox (As it is very computationally expensive)

→ Initial value of θ : For gradient descent & advanced optimization method, we need to set some initial theta value

→ Zero initialization : Initializing all theta weights to 0 doesn't work with neural networks. When we backpropagate, all nodes will update to the same value repeatedly resulting in same values of $\frac{\partial J(\theta)}{\partial \theta_{ij}}$ for each iteration of gradient descent.

And if the value of the weights does not change are same, that means over all hidden layers are computing the same exact feature. which will make this prevents our NN is pushing out complex hypothesis.

→ Random initialization : Symmetry breaking :

* Initialize each θ_{ij}^l to a random value in $[-\epsilon_{init}, \epsilon_{init}]$

R.P. $[-\epsilon_{init}, \epsilon_{init}]$ (i.e. $-\epsilon_{init} \leq \theta_{ij}^l \leq \epsilon_{init}$)

* Typically

$$\theta_{ij}^l = 2\epsilon_{init} \times \text{rand}(s_{l+1}, s_l + 1) - \epsilon_{init}$$

Max value = 1, Min = 0

The above expression keeps in check that θ_{ij}^l lies in the range $[-\epsilon_{init}, \epsilon_{init}]$

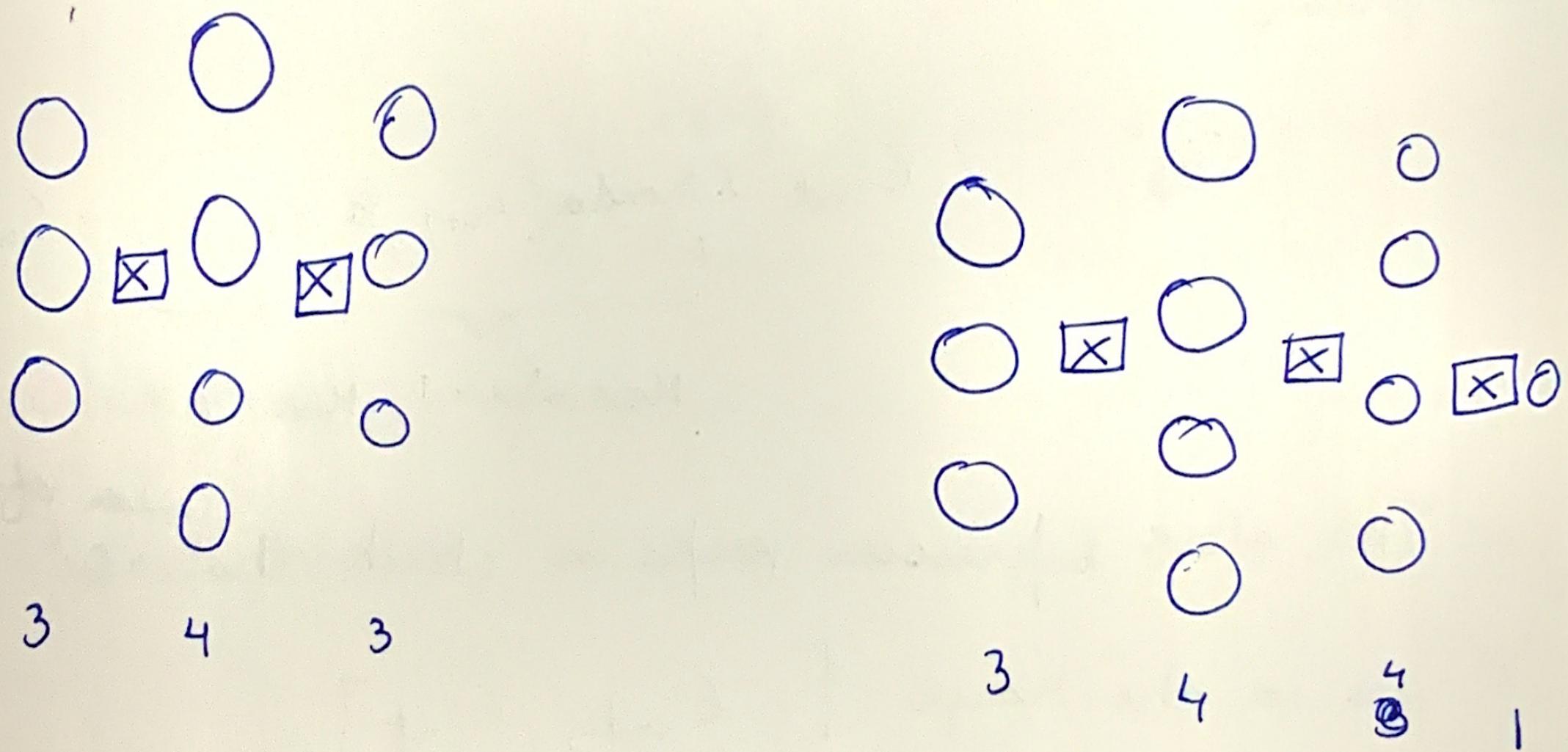
* Typically, we take ϵ as a small value like 0.1.

Because, if the value is too large, then our sigmoid activation keep putting out near to 1 values for each output (because of how large our weighted sum becomes).

Putting it together

→ Brainining a neural network

→ Pick a network architecture (Connectivity pattern b/w neurons)



→ No. of input units = dimensions of features x_i

→ No. of output units = no. of classes

→ No. of hidden units per layer : • Usually more the better. (Can increase computational complexity)

→ Rule: If more than 1 hidden layer, then assign the same no. of units in every layer

\Rightarrow Training a neural network

- 1) Randomly initialize the weights.
- 2) Implement forward propagation to get $h_\theta(x^i)$ for any x^i out of our training samples.
- 3) Implement code to compute cost function $J(\theta)$
- 4) Implement backprop to compute partial derivatives $\frac{\partial J(\theta)}{\partial \theta_{ij}}$

\rightarrow When we perform forward & backpropagation, we loop on every training example:

```
for i = 1 : m {
```

 Forward propagation on (x^i, y^i) // Get ~~b^a~~^a term

 Back propagation on (x^i, y^i) // Get δ terms

 Compute $\Delta := \Delta^l + \delta^{l+1} (a^l)^T$

```
}
```

5) Use gradient checking to compare $\frac{\partial J(\theta)}{\partial \theta_{ij}}$

Computed using backpropagation vs. using numerical estimate of gradient of $J(\theta)$.

Then disable gradient checking

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

* In neural nets, $J(\theta)$ is a non convex funⁿs thus we can end up in a local min^m minimum instead but it works even that might work very well for our problem.