


Unrolling Parameters:

May 19, 2021

* We need to "unroll" the Θ and D matrices to pass into MATLAB functions, since they only accept vectors.

e.g. $s_1 = 10, s_2 = 10, s_3 = 1$

$$\Theta^{(1)} \in [10 \times 1], \Theta^{(2)} \in [10 \times 1], \Theta^{(3)} \in [1 \times 1]$$

$$D^{(1)} \in [10 \times 1], D^{(2)} \in [10 \times 1], D^{(3)} \in [1 \times 1]$$

$$\rightarrow \text{ThetaVec} = [\text{Theta1}(:); \text{Theta2}(:); \text{Theta3}(:)];$$
$$DVec = [D1(:); D2(:); D3(:)];$$

Or:

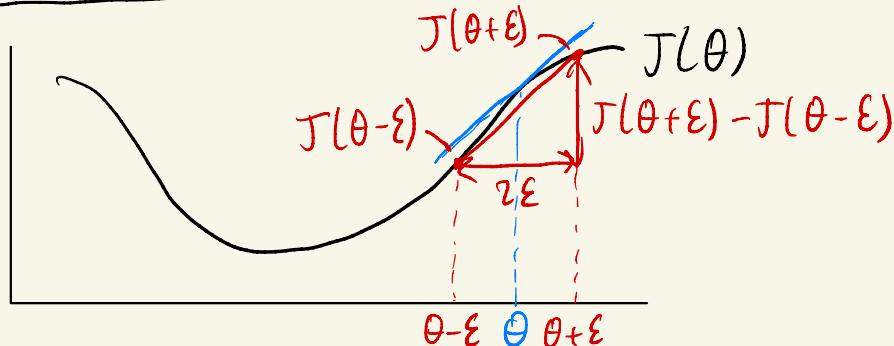
$$\text{Theta1} = \text{reshape}(\text{thetaVec}(1:110), 10, 11);$$

$$\text{Theta2} = \text{reshape}(\text{thetaVec}(111:220), 10, 11);$$

$$\text{Theta3} = \text{reshape}(\text{thetaVec}(221:231), 10, 11);$$

Numerical Estimation of Gradients:

e.g.



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

For Parameter Vector θ :

$\theta \in \mathbb{R}^n$ (unrolled)

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

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

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

$$\vdots$$

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \dots, \theta_{n-1}, \theta_n + \epsilon) - J(\theta_1, \theta_2, \dots, \theta_{n-1}, \theta_n - \epsilon)}{2\epsilon}$$

Finally, check that $\frac{\partial}{\partial \theta} J(\theta) \approx DVec$.

* Once you check backprop works, stop computing gradApprox since it is very slow.

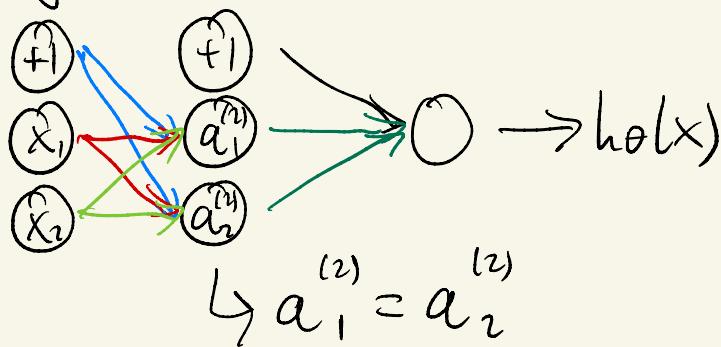
Random Initialization!

What do we initialize θ as?

↳ Can't set θ to all zeros (i.e. $\theta_{ij}^{(l)} \neq 0$ for all i, j, l)

↳ If they are, after each update, parameters corresponding to inputs going into all of the hidden units are identical.

e.g.



* Solve the problem of symmetric weights w/ random initialization \rightarrow symmetry breaking.

Symmetry Breaking

Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$.

e.g. \rightarrow all values of 10×11 matrix are $[0, 1]$

$$\Theta_{\text{theta1}} = \text{rand}(10, 11) \cdot (2 \cdot \text{INIT_EPSILON}) - \text{INIT_EPSILON}$$

$$\Theta_{\text{theta2}} = \text{rand}(1, 11) \cdot (2 \cdot \text{INIT_EPSILON}) - \text{INIT_EPSILON}$$

Putting It All Together:

1. Pick a neural network architecture.

\rightarrow # input units = dimension of features $x^{(i)}$

\rightarrow # output units = # of classes

\rightarrow 1 hidden layer, or if multiple, have same # of units in every hidden layer (usually more is better)

2. Randomly initialize weights

3. Implement forward prop to get $h_{\theta}(x^{(i)})$ for any $x^{(i)}$

4. Implement code to compute cost function $J(\theta)$

5. Implement backprop to compute partial derivatives

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

5. Use gradient checking to compare backprop computed gradient w/ numerical estimate.
→ Then disable gradient checking.
6. Use gradient descent w/ backprop to minimize $J(\Theta)$ as a function of parameters Θ .