

Cost Function:

For a neural Network:

L = total number of layers in the network

S_L = number of units in layer L (input)
(Not counting bias Unit)

K = Number of output units / classes

$$J(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[y_k^{(i)} \log((h_{\Theta}(x^{(i)}))_k) + (1 - y_k^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^m \sum_{j=1}^{S_{l+1}} (\Theta_{j,i}^{(l)})^2$$

Backpropagation Algorithm:

Given Training set = $\{(x^1, y^1), \dots, (x^m, y^m)\}$

Set, $\Delta_{i,j}^{(L)} = 0 \rightarrow$ Error Matrix

For training example $t=1$ to m :

\rightarrow set $a^{(1)} := x^{(t)}$

\rightarrow Perform Forward propagation a^l for $l=1, 2, \dots, L$

\rightarrow Using y^t , compute $\delta^{(L)} = a^L - y^t$

layer number
 \uparrow

→ Compute $\delta^{L-1}, \dots, \delta^2$ using,

$$\delta^{(l)} = ((\Theta^{(l)})^T \delta^{(l+1)}) \cdot \underbrace{a^{(l)} \cdot (1 - a^{(l)})}_{g'(z^{(l)})}$$

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

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

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

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$$

Unroll:

$$\text{thetavec} = [\text{Theta1}(:); \text{Theta2}(:); \text{Theta}(:);]$$

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

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

$$\rightarrow \text{Theta3} = \text{reshape}(\text{thetavec}(221:231), 1, 11)$$

Gradient Checking (Numerically):

```
epsilon = 1e-4;
```

for $i = 1 : n$,

$$\text{thetaplus} = \text{theta};$$
$$\text{thetaPlus}(i) += \text{epsilon};$$

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

$$\text{theta minus}(i) = \text{epsilon};$$
$$\text{grad Approx}(i) = \frac{(J(\text{theta plus}) - J(\text{theta minus}))}{(2 * \text{epsilon})};$$

end;

check grad Approx \approx DeltaVector

Random Initialization:

Initializing all thetas weight to ZERO does not work with neural network. Because in backpropagation, all nodes will be updated with same value repeatedly.

To break symmetry :

$$\text{Theta1} = \text{rand}(10, 11) * (2 * \text{EPS}) - \text{EPS};$$

$$\text{Theta2} = \text{rand}(1, 11) * (2 * \text{EPS}) - \text{EPS};$$

Big Picture of Neural Network:

Pick a network architecture.

- # Input Units = dimension of features $x^{(i)}$
- # Output Units = number of classes
- # hidden units per layer = the more the better
- Default 1 hidden layer, but again the more the better. [Some # of hidden units in all layer]

Training a Neural Network:

- Randomly Initialized the weights (Θ)
- Implement forward propagation to get $h_{\theta}(x^{(i)})$ for any $x^{(i)}$
- Implement the cost Function
- Implement back propagation to compute partial derivative

- Use gradient checking to confirm that back propagation works, then disable it.
- Use gradient descent or a built-in optimization function to minimize the cost function with weights in θ .