

WEEK-5

⇒ L = total no of layers in network

⇒ K = no of output units.

⇒ S_L = no of units in layer L (not counting bias unit)

⇒ COST FUNCTION → For Neural Networks

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

→ we don't sum over bias unit.

⇒ The double sum simply adds the logistic regression costs calculated for each cell in the output layer.

⇒ The triple sum adds up the squares of all the individual θ s in the entire network.

⇒ $\delta_j^{(l)}$ = 'error' of node j in layer l .

↳ $\delta_j^{(4)} = a_j^{(4)} - y_j$ → Back

Propagation
(way of finding $\frac{\partial}{\partial \theta} J(\theta)$)

$$\Rightarrow \delta^3 = (\theta^{(3)})^T \delta^4 * g'(z^{(3)})$$

$\hookrightarrow a^{(3)} * (1 - a^{(3)})$
 $\hookrightarrow \text{derivative of } g(z^{(3)})$

\Rightarrow So, using back propagation, we can compute gradient and find value of θ .

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)} \quad \left(\begin{array}{l} \text{ignoring } \lambda, \\ \text{if } \lambda = 0 \end{array} \right)$$

Back propagation Algorithm. (BP)

Step 1 \Rightarrow Training set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

S-2 \Rightarrow set $\Delta_{ij}^{(l)} = 0$, for all l, i, j (done for every $x^{(i)}, y^{(i)}$)

S-3 \Rightarrow For $i = 1$ to m :-
 set $a^{(1)} = x^{(i)}$

\rightarrow Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \dots, L$ (last layer)

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

\rightarrow using BP, compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^2$

$\rightarrow \Delta_{ij}^{(l)} = \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ ($\delta^{(1)}$ is not calculated)

⇒ Vectorized ⇒ $\Delta^{(l)} = A^{(l)} + \delta^{(l+1)} \cdot (a^{(l)})^T$

S-4 ⇒ $D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)} + \frac{\lambda}{m} \theta_{ij}^{(l)} \text{ if } j \neq 0$

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

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

Vectorized ⇒ $\text{theta_grad} = (1/m) * \text{del } J / \text{del } \theta_{a1j}$

⇒ Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i) \text{ (for } j \geq 0 \text{)}$

where $\text{cost}(i) = y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))$

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

$\delta_2^{(2)} = \theta_{22}^{(2)} \delta_2^{(3)}$

⇒ for example, if $s_1 = 10, s_2 = 10, s_3 = 1$

⇒ $\theta^1 \Rightarrow \mathbb{R}^{10 \times 11}, \theta^2 \Rightarrow \mathbb{R}^{10 \times 11}, \theta^{(3)} = 1 \times 11$

⇒ $D^{(1)} \Rightarrow \mathbb{R}^{10 \times 11}, D^{(2)} \Rightarrow \mathbb{R}^{10 \times 11}, D^{(3)} = 1 \times 11$

⇒ $\text{theta_vec} = [\text{theta1}(:); \text{theta2}(:); \text{theta3}(:)];$
converts these vectors in one long vector

⇒ to redo it & get theta1

Also called
Gradient
Checking

theta1 = reshape(thetaVec (1:110), 10, 11);

⇒ to use these in optimization f's like
fminunc() → we ~~wrap~~ roll them into one
big long vector.

⇒ gradApprox = $\frac{J(\text{theta} + \text{epsilon}) - J(\text{theta} - \text{epsilon})}{2 \times \text{epsilon}}$
Way to check BP

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

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

⇒ $\frac{\partial}{\partial \theta_i} J(\theta)$ for every θ .

$\epsilon = 10^{-4}$

θ

> (Gradient Checking)
Code

⇒ for $i = 1:n$,
 thetaPlus = theta;
 thetaPlus(i) = thetaPlus(i) + epsilon;
 thetaMinus = theta;
 thetaMinus(i) = thetaMinus(i) - epsilon;
 gradApprox(i) = $\frac{J(\text{thetaPlus}) - J(\text{thetaMinus})}{(2 \times \text{epsilon})}$;
end;

⇒ if (gradApprox ≈ 0Vec) → then BP is rightly
 used

⇒ Gradient checking is very slow as compared to BP. (so only 1/2 once)

⇒ If $\theta_{ij}^{(l)} = 0$ for all i, j, l

$$\therefore a_4^{(2)} = a_2^{(2)}, \quad \delta_1^{(2)} = \delta_2^{(2)}$$

$$\frac{\partial J(\theta)}{\partial \theta_{01}^{(1)}} = \frac{\partial J(\theta)}{\partial \theta_{02}^{(1)}}, \quad \theta_{01}^{(1)} = \theta_{02}^{(1)}$$

⇒ So we can't initialize the theta (parameters) with 0 (as in logistic regression), therefore we use random values (Symmetry Breaking)

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

Code

→ To create a random 10×11 matrix (values between -1 and 1)

$$\text{theta1} = \text{rand}(10, 11) * (2 * \text{init_epsilon}) - \text{init_epsilon}$$

$$\text{theta2} = \text{rand}(1, 11) * (2 * \text{init_epsilon}) - \text{init_epsilon} \quad \left. \begin{array}{l} \text{To keep the} \\ \text{value in} \\ [-\epsilon, \epsilon] \text{ range} \end{array} \right\}$$

(this epsilon is unrelated to epsilon in gradient checking)

- \Rightarrow No. of input units \Rightarrow Dimension of features (x^i)
- \Rightarrow No. of output units \Rightarrow Number of classes.
- \Rightarrow For hidden layers ≥ 1 or if > 1 hidden layers, should have same no. of hidden units or nodes in every layer.

\Rightarrow TRAINING A Neural Network

- 1) = Randomly initialize weights
- 2) = Implement forward propagation to get $h_\theta(x^{(i)})$ for any $x^{(i)}$
- 3) = Implement code to compute cost $J \Rightarrow J(\theta)$
- 4) = Implement backprop to compute partial derivatives $\rightarrow \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$

for $i = 1:m$

\rightarrow Perform forward P & BP using

example (x^i, y^i)

\rightarrow get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l = 2, \dots, L$.

- 5) = Use gradient checking to compare $\frac{\partial}{\partial \theta} J(\theta)$ computed using BP v.s. using numerical estimate of gradient of $J(\theta)$. Disable GC

6) = Use gradient descent or advanced optimization method with BP to try to minimize $J(\theta)$ as a fn of θ .

\Rightarrow For NN, $J(\theta)$ is non-convex.