

Neural Networks Cost Function

Binary classification

Multi-class classification (K classes)

1 output unit

K output units

$$h_{\theta}(x) \in \mathbb{R}$$

$$h_{\theta}(x) \in \mathbb{R}^K$$

$$S_L = 1 \quad K = 1$$

$$S_L = K \quad (K \geq 3)$$

Cost Function

$$h_{\theta}(x) \in \mathbb{R}^K \quad (h_{\theta}(x))_i = i^{\text{th}} \text{ output} \quad y_K = \begin{bmatrix} 0 \\ 0 \\ \vdots \end{bmatrix} \dots$$

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

Backpropagation Algorithm (for derivatives)

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

$$\min_{\theta} J(\theta)$$

Need code to compute

$$J(\theta)$$

$$\frac{\partial}{\partial \theta_{ij}} J(\theta)$$

Intuition: $\delta_j^{(l)}$ = "error of node j in layer l "

$$\delta_j^{(4)} = \boxed{a_j^{(4)}} - y_j \quad (h_{\theta}(x))_j \rightarrow \text{vectorize} \quad \delta^{(4)} = a^{(4)} - y$$

$$\delta^{(3)} = (\theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)}) \quad a^{(3)} \cdot (1 - a^{(3)})$$

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} \cdot g'(z^{(2)}) \quad a^{(2)} \cdot (1 - a^{(2)})$$

(No $\delta^{(1)}$)

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

Set $\Delta_{ij}^{(L)} = 0$ (for l, i, j) (next to compute $\frac{\partial}{\partial \theta_{ij}^{(L)}} J(\theta)$)

For $i=1$ to m training example $(x^{(i)}, y^{(i)})$

Set $a^{(1)} = x^{(i)}$

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

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

Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ ~~$\delta^{(1)}$~~

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)} \rightarrow \text{vectorized} \quad \Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

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

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

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

Backpropagation Intuition

$\ominus_{ij}^{(l)} \leftarrow$ layer
to κ from

Implementation note: Unrolling Parameters

Gradient Checking

$$\frac{\partial}{\partial \theta} J(\theta) \approx \frac{J(\theta+\epsilon) - J(\theta-\epsilon)}{2\epsilon} \leftarrow \begin{matrix} \text{2-sided} \\ \text{diff} \end{matrix} \quad \epsilon = 10^{-4}$$

$$\frac{J(\theta+\epsilon) - J(\theta)}{\epsilon} \leftarrow \begin{matrix} \text{1-sided} \\ \text{diff} \end{matrix}$$

$$\Theta \in \mathbb{R}^n$$

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

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



Random Initialization

need initial value for Θ

Zero initialization, makes $a_1^{(2)} = a_2^{(2)}$ & $S_1^{(2)} = S_2^{(2)}$ (TAID)

(computes all the same feature)

random within $[-\epsilon, \epsilon]$

Putting it Together

input units: Dimension of features $x^{(i)}$

output units: number of classes

Reasonable default: 1 hidden layer, if $R \geq 1$ hidden layer, \rightarrow same # of units in each

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 function $J(\Theta)$
4. Implement backprop to compute partial derivatives $\frac{\partial}{\partial \Theta_{j,k}^{(l)}} J(\Theta)$
5. Use gradient checking to compare $\frac{\partial}{\partial \Theta_{j,k}^{(l)}} J(\Theta)$ (computed using backprop) vs. using estimate of gradient
- Disable gradient checking code
6. Use gradient descent or advanced optimization method w/ backprop to minimize $J(\Theta)$

Autonomous Driving