

Neural Network classification problem

training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

L : total no. of layers in the network

s_l : no. of units (exclude bias unit) in layer l

binary classification: $y = 0 \text{ or } 1$ | output unit
 $\begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$

Multi-class classification (K classes) $\begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$ K output units

logistic regression

$$J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\Theta(x^{(i)}))_k + (1-y_k^{(i)}) \log(1-h_\Theta(x^{(i)}))_k \right] - \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^{s_i} \sum_{k=1}^{s_{i+1}} (h_\Theta^{(i)}(x^{(i)})_{jk})^2$$

$h_\Theta(x) \in \mathbb{R}^K$ $(h_\Theta(x))_i$: i^{th} output

Neural Networks: Backpropagation algorithm

Gradient Computation: need $J(\Theta)$ $\frac{\partial J(\Theta)}{\partial \Theta_{ij}^k}$

e.g. 1 training example, 4 layers, forward propagation:

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

$$z^{(2)} = \Theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$

$$z^{(3)} = \Theta^{(2)} a^{(2)} \quad a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0^{(3)})$$

$$z^{(4)} = \Theta^{(3)} a^{(3)} \quad a^{(4)} = h_\Theta(x) = g(z^{(4)})$$

Backpropagation algorithm

intuition: $\delta_j^{(l)} = \text{"error" of node } j \text{ on layer } l$

$$\text{e.g. for layer 4: } \delta_j^{(4)} = a_j^{(4)} - y_j$$

$$\delta_j^{(3)} = (\Theta^{(3)})^T \delta_i^{(4)} * g'(z^{(3)}) \quad a_i^{(3)} * (1 - a_i^{(3)})$$

$$\delta_j^{(2)} = (\Theta^{(2)})^T \delta_i^{(3)} * g'(z^{(2)}) \quad a_i^{(2)} * (1 - a_i^{(2)})$$

no $\delta^{(1)}$ term

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

Backpropagation Algorithm

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

set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j) (use to compute $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\Theta)$)

for $i=1$ to m :

↳ set $a^{(1)} = X^{(i)}$ (input layer)

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

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

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

↳ $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

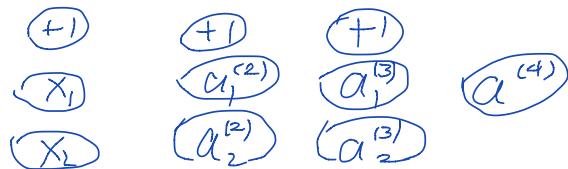
$\rightarrow D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \quad \text{if } j \neq 0$

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

\rightarrow use for gradient descent or other algorithms

Backpropagation intuition

forward propagation



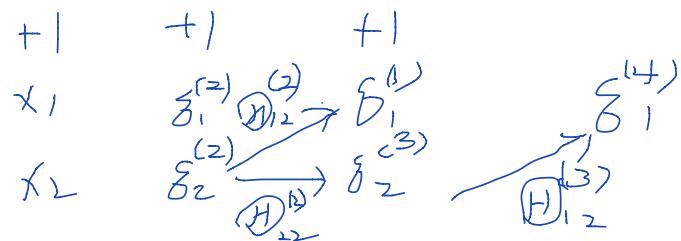
$$z_1^{(3)} = \Theta_{10}^{(4)} x_1 + \Theta_{11}^{(4)} a_1^{(2)} + \Theta_{12}^{(4)} a_2^{(2)}$$

backpropagation

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

$\delta_j^{(l)}$ = "error" of cost for $a_j^{(l)}$ (unit j in layer l)

$$\text{formally: } \delta_j^{(l)} = \frac{\partial}{\partial z_j} \text{cost}(1)$$



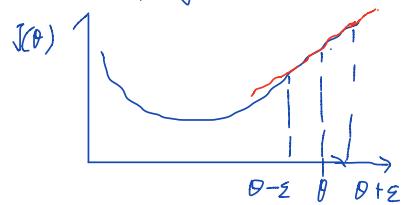
$$\delta^{(4)} = y^{(4)} - a_1^{(4)}$$

$$\delta_2^{(3)} = \Theta_{12}^{(3)} \cdot \delta_1^{(4)}$$

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

Gradient Checking

Numerical estimation of gradients



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

$$\epsilon \sim 10^{-4}$$

Parameter vector θ

$\theta \in \mathbb{R}^n$ (unrolled version of $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$)

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

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

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \dots \quad \dots \quad \dots$$

⋮

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \dots \quad \dots \quad \dots$$

check that above approx $\approx \frac{\partial}{\partial \theta_j} J(\theta)$ from back propagation

Implementation Note

- 1) Implement backprop. to compute DVec (unrolled $\vec{J}^{(1)}, \vec{J}^{(2)}, \dots$)
- 2) Implement numerical gradient check to compute gradApprox
- 3) Make sure they give similar values
- 4) Turn off gradient checking, using backprop code for learning.

Random Initialization

Initial value of θ

- ★. Initialize all parameters to 0 when training a

Neural network is NOT OK !

(after each update, parameters corresponding to inputs going into each of two hidden units are identical
this prevents your neural network from learning anything interesting)

Random Initialization: Symmetry breaking

initialize each Θ_{ij}^w to a random value in $[-\varepsilon, \varepsilon]$

Put it together.

1) Pick a network architecture

No. of input units: Dimension of features $X^{(1)}$

No. of output units: number of classes

No. of hidden layers:

reasonable default: 1 hidden layer

or if > 1 hidden layer, have same number of hidden units in each layer (Usually the more -the better)

2) Randomly initialize weights Θ

3) Implement forward propagation to get $h_\Theta(x^{(1)})$ for any $x^{(1)}$

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

5) Implement backprop to compute partial derivatives $\frac{\partial}{\partial \Theta_k} J(\Theta)$

for $i = 1 : m$ $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$

perform forward prop and backprop using example

$(x^{(1)}, y^{(1)})$, get activations $a^{(1)}$ and delta $\delta^{(1)}$

for $l = 2, \dots, L$

- 6) Use gradient checking to compare $\frac{\partial^2 J(\theta)}{\partial \theta_j \partial \theta_k}$ from backprop and numerical estimation, then disable the gradient checking
- 7) Use gradient descent or adv. method with backprop to try to minimize $J(\theta)$ as a function of parameters θ