

Neural Network Learning

KAUSHIK ROY

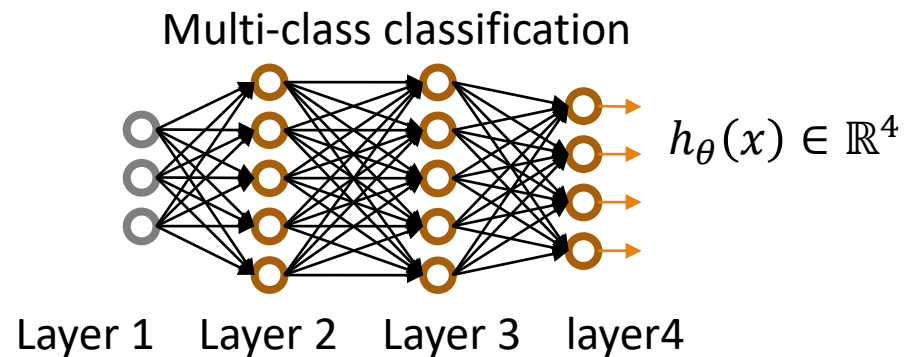
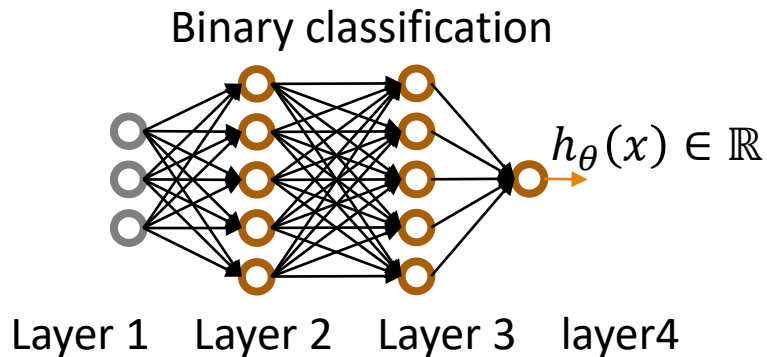
Neural Networks

➤ Classification

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

L = total no. of layers in the network

s_l = no. of units (not counting bias unit) in layer l



Each output neuron performs a binary classification task.
The number of output neurons = number of classes.

Neural Networks

➤ Cost function

Recall the logistic regression cost function:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[\log(h_{\theta}(x^{(i)})) * y^{(i)} + \log(1 - h_{\theta}(x^{(i)})) * (1 - y^{(i)}) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Neural network cost function:

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

In which, $h_{\theta}(x) \in \mathbb{R}^K$, and $(h_{\theta}(x))_k = k^{th}$ output, and $1 \leq k \leq K$ (outputs)

Neural Networks

➤ Gradient computation

Recall Neuron network cost function:

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

In which, $h_{\theta}(x) \in \mathbb{R}^K$, and $\left(h_{\theta}(x) \right)_k = k^{th}$ output, and $1 \leq k \leq K$

To minimize $J(\theta)$, we need to compute the gradient: $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$

Neural Networks

➤ Gradient computation

E.g. Assume a 4 layer network and one training example (x, y) :

In the forward path:

Recall forward propagation:

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

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

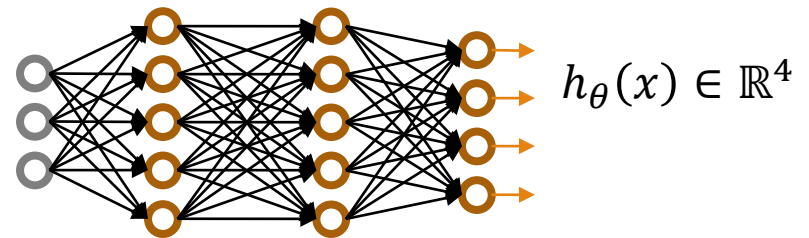
$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$

$$z^{(3)} = \theta^{(2)} a^{(2)}$$

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

$$z^{(4)} = \theta^{(3)} a^{(3)}$$

$$a^{(4)} = h_{\theta}(x) = g(z^{(4)})$$



Layer 1 Layer 2 Layer 3 layer4

Dimensions equal number of neurons in the layer

E.g.

$$a^{(1)} \in \mathbb{R}^3$$

$$a^{(2)} \in \mathbb{R}^5$$

$$a^{(3)} \in \mathbb{R}^5$$

$$a^{(4)} \in \mathbb{R}^4$$

Neural Networks

➤ Gradient computation

In the backward path:

Intuition:

$$\delta_j^{(l)} = \text{“error” of node } j \text{ in layer } l$$

For each output unit (layer L=4)

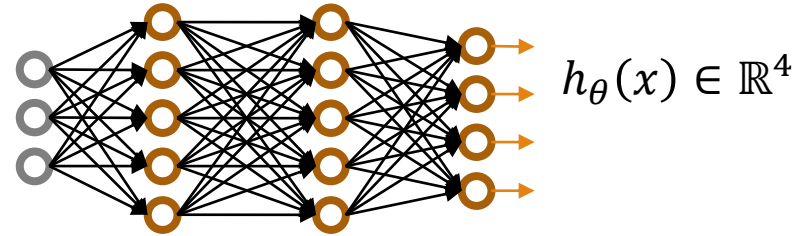
$$\delta^{(4)} = a^{(4)} - y = (h_{\theta}(x)) - y$$

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

$$\delta^{(2)} = (\theta^{(2)})^T \delta^{(3)} .* g'(z^{(2)})$$

No need to compute $\delta^{(1)}$ for input layer

One can show: $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)}$



Layer 1 Layer 2 Layer 3 layer4

Dimensions equal number of neurons in the layer

E.g.

$$\delta^{(4)} \in \mathbb{R}^4$$

$$\delta^{(3)} \in \mathbb{R}^5$$

$$\delta^{(2)} \in \mathbb{R}^5$$

$$g'(z^{(3)}) = a^{(3)} .* (1 - a^{(3)})$$

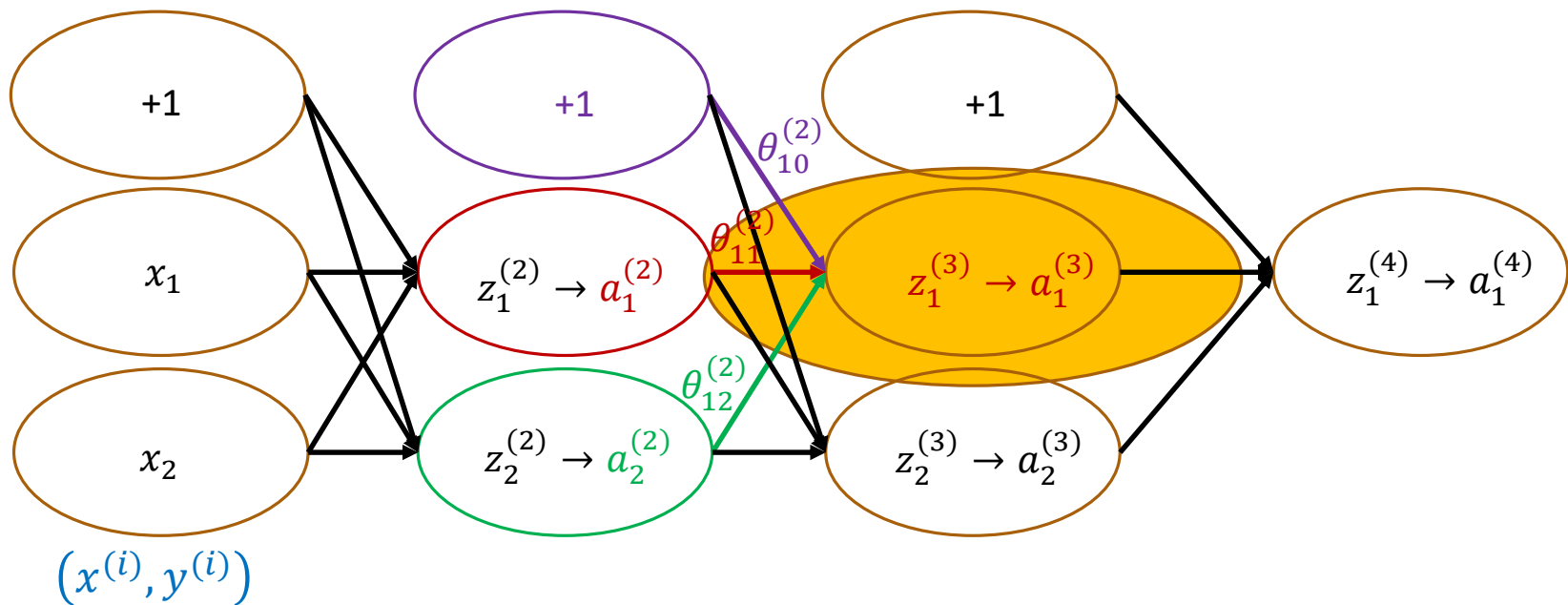
$$g'(z^{(2)}) = a^{(2)} .* (1 - a^{(2)})$$

Backpropagation

➤ Forward Propagation

Forward path:

In the forward path of backpropagation, input data propagates in the forward direction layer by layer until it reaches the output layer. The input to one layer is the output of the previous layer.



E.g. Outputs from layer 2 are used to compute $z_1^{(3)}$ in layer 3:

$$z_1^{(3)} = \theta_{10}^{(2)} \times 1 + \theta_{11}^{(2)} \times a_1^{(2)} + \theta_{12}^{(2)} \times a_2^{(2)}$$

Backpropagation

➤ Backpropagation

What is backpropagation doing?

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

Focusing on a single example $x^{(i)}, y^{(i)}$, the case of 1 output unit, and ignoring regularization ($\lambda = 0$),

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

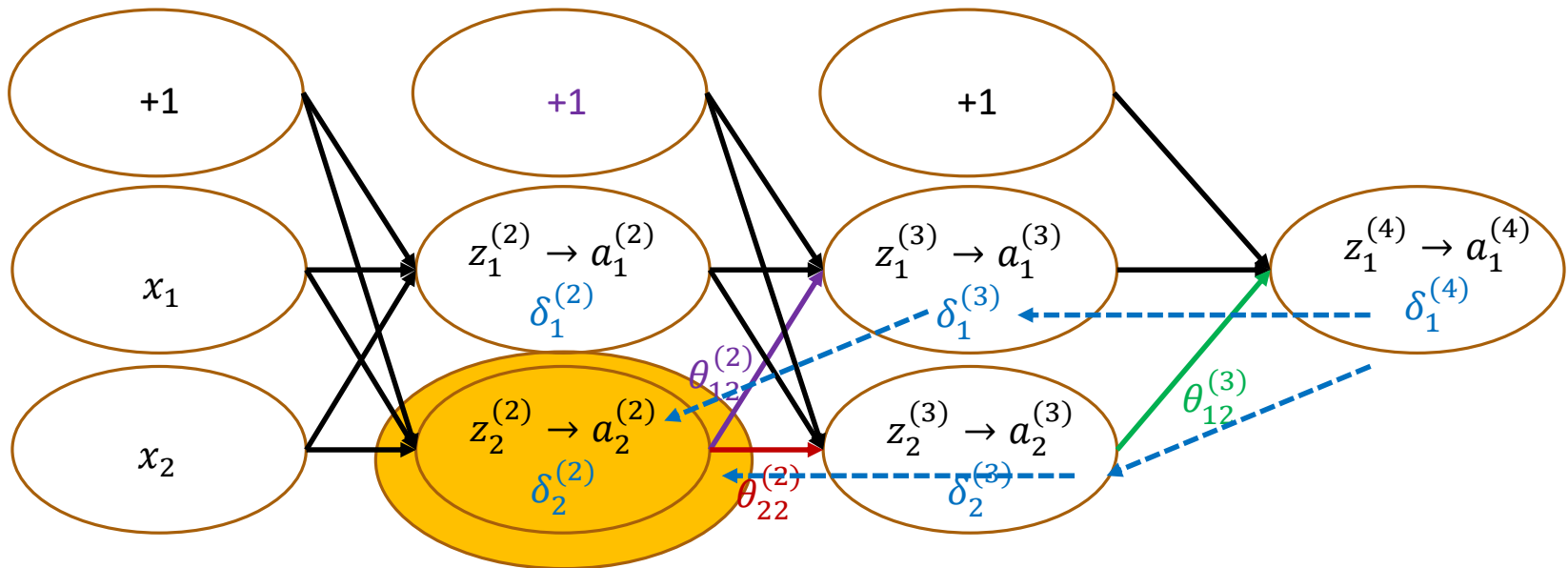
(Think of $cost(i) \approx (h_{\theta}(x^{(i)}) - y^{(i)})^2$)

How well is the network doing on example i ?

Backpropagation

➤ Backpropagation

Backward path:

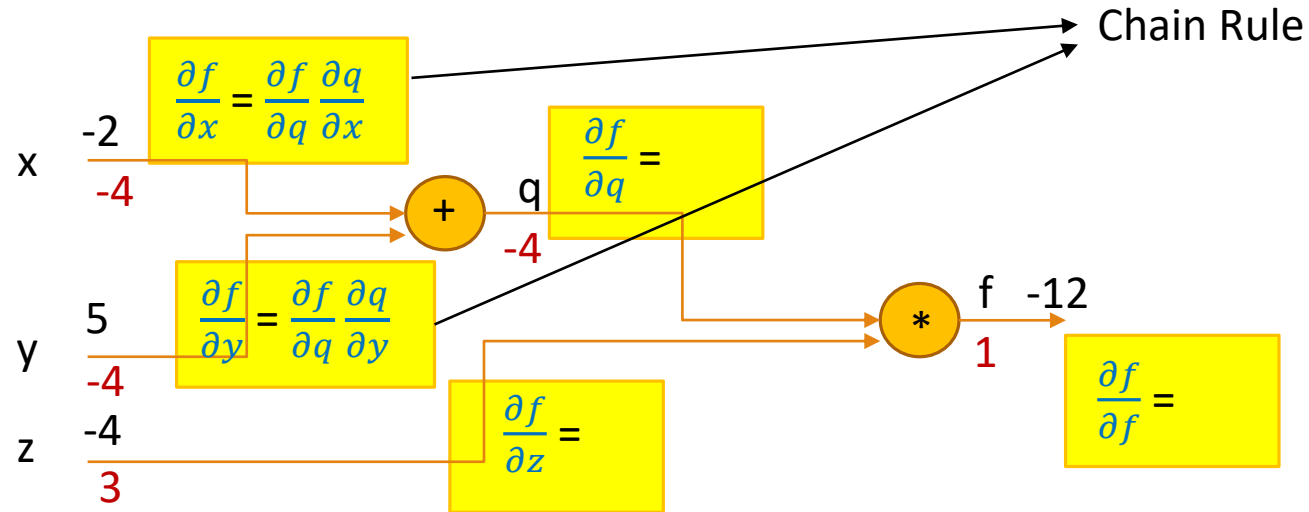


$\delta_j^{(l)}$ = "error" of cost for $a_j^{(l)}$ (unit j in layer l)
 Formally, $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} \text{cost}(i)$ (for $j \geq 0$), where
 $\text{cost}(i) = y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log h_{\theta}(x^{(i)})$

E.g. $\delta_1^{(4)} = y^{(i)} - a_1^{(4)}$
 $\delta_2^{(3)} = \theta_{12}^{(3)} \delta_1^{(4)}$
 $\delta_2^{(2)} = \theta_{12}^{(2)} \delta_1^{(3)} + \theta_{22}^{(2)} \delta_2^{(3)}$

Example Backpropagation

➤ $f(x,y,z) = (x + y)z$



$$q = (x + y) \quad \frac{\partial q}{\partial x} = 1 \quad \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z \quad \frac{\partial f}{\partial z} = q$$

Example:

➤ $f(w,x) = \text{sigmoid}$

Backpropagation

➤ Backpropagation algorithm

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

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

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

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

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

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

Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ (using BP; no need to compute $\delta^{(1)}$)

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$ (possible to vectorize: $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^\top$)

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

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

If can be shown that $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = D_{ij}^{(l)}$