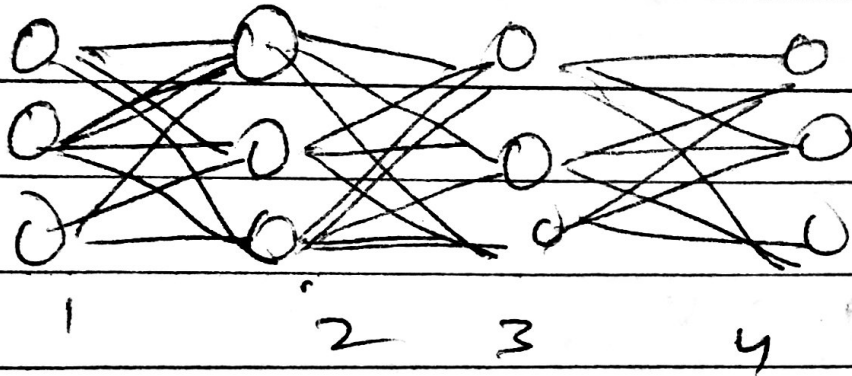


5

Neural Network (classification)



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

L = total no. of layers in network ($= 4$)

S_l = no. of units (excluding bias) in layer l .

$$S_1 = 3, S_2 = 3, S_4 = 3$$

Binary classification
 $y = 0$ or 1
1 output units

Multi classes (k classes)

k output units

* Don't regularize bias term

Logistic regression:

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

Regularization term

Neural Network:

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

Regularization for $\theta_{ji}^{(k)}$

Cost function

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

→ we are finding theta to minimize $J(\theta)$

compute

$$J(\theta) \quad \& \quad \frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$$

Remember

[forward propagation]

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

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

$$a^{(2)} = g(z^{(2)}) \dots [\text{add } a_0^{(2)}]$$

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

$$a^{(3)} = g(z^{(3)}) \dots [\text{add } a_0^{(2)}]$$

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

$$a^0 = h_{\theta}(x) = g(z^{(4)})$$

↑ 4 layers only

for derivatives, we use back propagation

we calculate

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

$$\delta_j^{(4)} = \boxed{a_j^{(4)}} - y_j \quad \leftarrow \begin{array}{l} \text{actual} \\ \text{value} \end{array}$$

$\nwarrow (h_\theta(x))_j$

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

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

$g' = \text{derivative}$

No $\delta^{(1)}$ term \rightarrow it has input

so you don't need to calculate (δ)

So now,

$$\frac{\partial}{\partial \theta^{(l)}} J(\theta) = a_j^{(l)} \delta_j^{(l+1)}$$

→ if we ignore the λ terms
i.e., if we ignore
the regularization term

What we do,

$$\text{set } \Delta_{ij}^{(L)} = 0 \dots (\text{not } \delta)$$

$$\uparrow (\text{and so compute } \frac{\partial}{\partial \theta_{ij}^{(L)}} J(\theta))$$

for $i = 1$ to m

$$\text{set } a^{(1)} = x^{(i)}$$

Now, compute forward propagation

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

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

compute $\delta^{(L+1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

Back prop

$$\Delta_{ij}^{(L)} = \Delta_{ij}^{(L)} + a_j^{(L)} \delta_i^{(L+1)}$$

↑

$$\text{vectorized } \Delta^L = \Delta^L + \delta^{(L+1)} (a^{(L)})^T$$

end for

After for loop

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

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

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

So now, you can use these in some gradient descent or advanced optimization algorithm