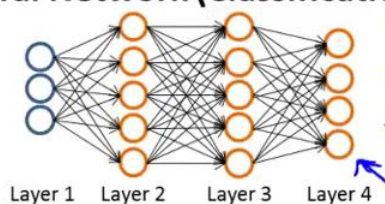


Neural Network (Classification)



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

$$\rightarrow L = \text{total no. of layers in network} \quad \underline{L=4}$$

$$\rightarrow s_l = \text{no. of units (not counting bias unit) in layer } l \quad s_1=3, s_2=5, s_3=5, s_4=4$$

Binary classification

$$y = 0 \text{ or } 1 \leftarrow$$

1 output unit \leftarrow

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

$$s_1 = 1, \quad K=1$$



Multi-class classification (K classes)

$$y \in \mathbb{R}^K \quad \text{E.g.} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \leftarrow$$

pedestrian car motorcycle truck

K output units

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

$$s_1 = K \quad (K \geq 3)$$

Cost function

Logistic regression:

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

Neural network:

$$\rightarrow h_{\Theta}(x) \in \mathbb{R}^K \quad (h_{\Theta}(x))_i = i^{th} \text{ output}$$

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

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2$$

$$\begin{matrix} \textcircled{1}^{(1)}_{ji} \\ \textcircled{1}^{(2)}_{i0} x_0 + \textcircled{1}^{(2)}_{i1} x_1 + \dots \end{matrix}$$

$$\sum_{j=1}^{s_l} \textcircled{1}^{(l)}_{ji} x_j = y_k \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Andrew

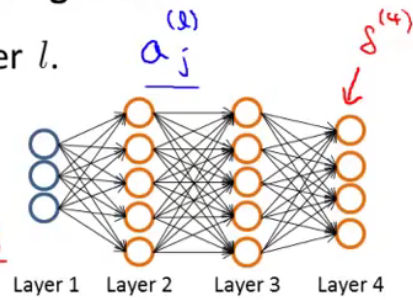
反向传播算法 Backpropagation

Gradient computation: Backpropagation algorithm

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

For each output unit (layer $L = 4$)

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



$$\delta_j^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)})$$

$$\delta_j^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \cdot g'(z^{(2)})$$

$$a^{(3)} \cdot (1 - a^{(3)})$$

$$a^{(2)} \cdot (1 - a^{(2)})$$

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 $m \leftarrow (x^{(i)}, y^{(i)})$.

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

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

$$\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$$

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

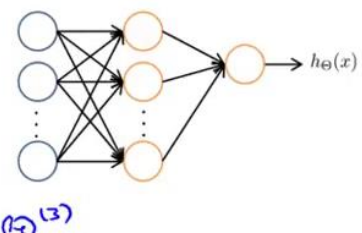
$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \text{ if } j = 0$$

Example

$$s_1 = 10, s_2 = 10, s_3 = 1$$

$$\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$$

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



> `thetaVec = [Theta1(:); Theta2(:); Theta3(:)];`

> `DVec = [D1(:); D2(:); D3(:)];`

`Theta1 = reshape(thetaVec(1:110), 10, 11);`

> `Theta2 = reshape(thetaVec(111:220), 10, 11);`

> `Theta3 = reshape(thetaVec(221:231), 1, 11);`

Learning Algorithm

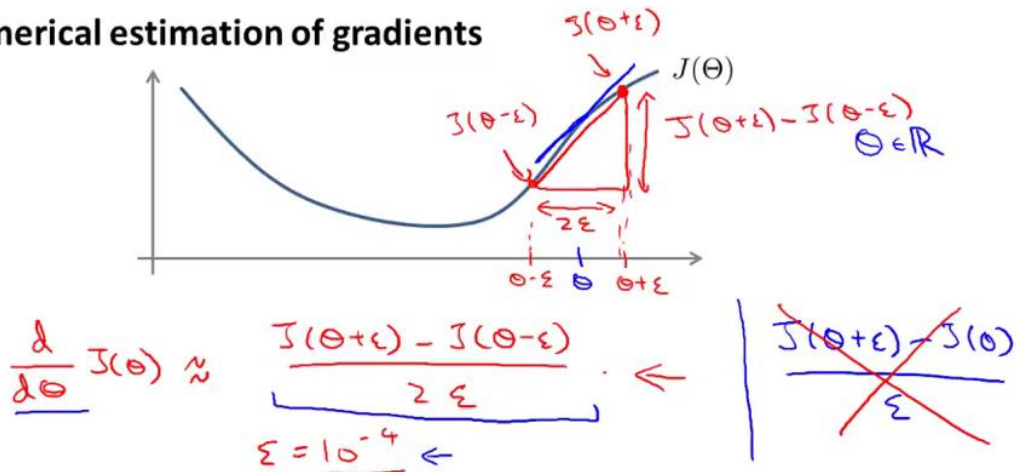
- Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.
- Unroll to get `initialTheta` to pass to
- `fminunc(@costFunction, initialTheta, options)`

`function [jval, gradientVec] = costFunction(thetaVec)`

- From `thetaVec`, get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ *reshape*
- Use forward prop/back prop to compute $D^{(1)}, D^{(2)}, D^{(3)}$ and $J(\Theta)$.
Unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get `gradientVec`.

梯度 checking

Numerical estimation of gradients



Implementation: `gradApprox = (J(theta + EPSILON) - J(theta - EPSILON)) / (2*EPSILON)`

for $i = 1:n$, ←

```

    thetaPlus = theta;
    thetaPlus(i) = thetaPlus(i) + EPSILON;
    thetaMinus = theta;
    thetaMinus(i) = thetaMinus(i) - EPSILON;
    gradApprox(i) = (J(thetaPlus) - J(thetaMinus)) / (2*EPSILON);

```

end;

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i + \epsilon \\ \vdots \\ \theta_n \end{bmatrix} \rightarrow \theta_i + \epsilon$$

$$\frac{2}{2\theta_i} J(\theta)$$

Check that `gradApprox` \approx `DVec`

From backprop.

Implementation Note:

- > - Implement backprop to compute DVec (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$).
- > - Implement numerical gradient check to compute gradApprox.
- > - Make sure they give similar values.
- > - Turn off gradient checking. Using backprop code for learning.

Important:

- > - Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of `costFunction(...)`) your code will be very slow.

随机初始化

Random initialization: Symmetry breaking

- > Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$ (i.e. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)

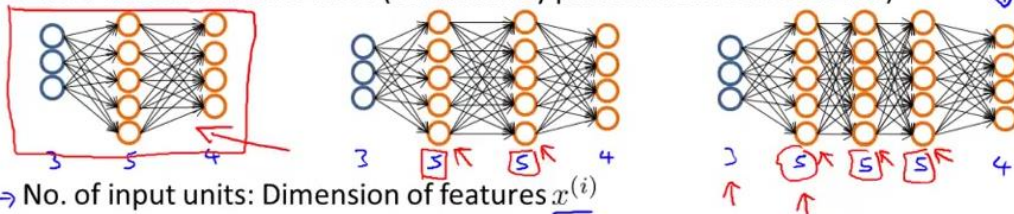
E.g.

→ `Theta1 = rand(10,11) * (2*INIT_EPSILON) - INIT_EPSILON;` $[-\epsilon, \epsilon]$

`Theta2 = rand(1,11) * (2*INIT_EPSILON) - INIT_EPSILON;`

Training a neural network

Pick a network architecture (connectivity pattern between neurons)



→ No. of input units: Dimension of features $x^{(i)}$

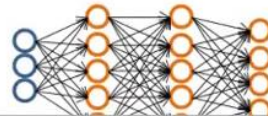
→ No. output units: Number of classes

Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

1 7 5 0 7 1 0 7

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 $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$
- for $i = 1:m$ $(x^{(1)}, y^{(1)}) \quad (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$
 - Perform forward propagation and backpropagation using example $(x^{(i)}, y^{(i)})$
(Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l = 2, \dots, L$).



Training a neural network

- 5. Use gradient checking to compare $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$ computed using backpropagation vs. using numerical estimate of gradient of $J(\Theta)$.
Then disable gradient checking code.
- 6. Use gradient descent or advanced optimization method with backpropagation to try to minimize $J(\Theta)$ as a function of parameters Θ