# ①. COST FUNCTION AND BACKPROPAGATION.

## 1.1. Cost Functn:

NN (Classificatn)

$L$ = Total # of layers in network.

$S_\ell$ = # of units in layer $\ell$ (excluding bias unit)

Binary classificatn: 1 output unit.   $h_\Theta(x) \in \mathbb{R}$

Multi-Class classfctn: K output units.   $h_\Theta(x) \in \mathbb{R}^K$

### Cost Func:

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

### NN:

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

## 1.2. Backpropagatn Algorithm:

$\min\limits_{\Theta} J(\Theta)$

Need code to compute : $\rightarrow J(\Theta)$

$\rightarrow -\frac{\partial}{\partial\Theta_{ij}^{(l)}} J(\Theta)$

### Gradient Computatn:

$(x,y) \leftarrow$ training example.

$a^{(1)} = x$

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

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

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

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

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

$a^{(4)} = g(z^{(4)}) = h_\Theta(x)$  (add $a_0^{(4)}$)

### Backpropagatn Algorithm:

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

For each output unit (layer L=4)

$\delta_j^{(4)} = a_j^{(4)} - y_j \longrightarrow \delta^{(4)} = a^{(4)} - y$

$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \ .* \ g'(z^{(3)}) \longrightarrow a^{(3)} .* (1-a^{(3)})$

$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \ .* \ g'(z^{(2)}) \longrightarrow a^{(2)} .* (1-a^{(2)})$

$\frac{\partial}{\partial\Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \cdot \delta_i^{(l+1)}$  (ignoring $\lambda$, if $\lambda=0$)

For the huge # of training set.

$$\Delta_{ij}^{(\ell)} = 0 \quad (\text{for } \forall i,j,\ell) \qquad (\text{used to compute } \frac{\partial}{\partial \Theta_{ij}^{(\ell)}} J(\theta))$$

for $i=1$ to $m$

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

    Perform forward propagatn to compute $a^{(\ell)}$ for $\ell = 2, 3, \ldots, L$

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

    compute $\delta^{(L-1)}, \delta^{(L-2)}, \ldots, \delta^{(2)} \quad \cancel{\delta^{(1)}}$

$$\Delta_{ij}^{(\ell)} = \Delta_{ij}^{(\ell)} + a_j^{(\ell)} \delta_i^{(\ell+1)} \quad \longrightarrow \quad \Delta^{(\ell)} = \Delta^{(\ell)} + \delta^{(\ell+1)} (a^{(\ell)})^T.$$

$$\left. \begin{aligned} D_{ij}^{(\ell)} &= \tfrac{1}{m} \Delta_{ij}^{(\ell)} + \lambda \Theta_{ij}^{(\ell)} \quad \text{if } j \neq 0 \\[8pt] D_{ij}^{(\ell)} &= \tfrac{1}{m} \Delta_{ij}^{(\ell)} \qquad\qquad \text{if } j = 0. \end{aligned} \right\} \quad \frac{\partial}{\partial \Theta_{ij}^{(\ell)}} J(\theta) = D_{ij}^{(\ell)}$$

end.

### 1.3. Backpropagatn Intuitn:

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

Backpropagatn: (running the feedforward algorithm, but doing it backwards.)

$$\text{cost}(i) = y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log h_\theta(x^{(i)}) \qquad \left[ \text{Think of cost}(i) \approx (h_\theta(x^{(i)}) - y^{(i)})^2 \right]$$

how well is the network doing on example $i$?

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

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

## ② BACKPROPAGATION IN PRACTICE.

### 2.1. Implementatn Note: Unrolling Paramtrs:

NN: $(L=4)$

$\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ — matrices (Theta 1, Theta 2, Theta 3)

$D^{(1)}, D^{(2)}, D^{(3)}$ — matrices (D1, D2, D3)

"Unroll" into vectors.

ex/ $S_1 = 10, \; S_2 = 10, \; S_3 = 1$

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

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

```
thetaVec = [Theta1(:); Theta2(:); Theta3(:)];
DVec = [D1(:); D2(:); D3(:)];
Theta1 = reshape(thetaVec(1:110), 10, 11);
```

## Learning Algorithm:

- have initial paramtrs $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$
- unroll to get initial Theta to pass to

$$fmm unc (@ costFunction, initialTheta, options).$$

function [jval, gradientVec] = costFunction (thetaVec)

  - from thetaVec, get $\theta^{(1)}, \theta^{(2)}, \theta^{(2)}$
  - 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.

### 2.2 Gradient Checking:

It may look $J(\theta)$ is decreasing, but you might end up with a NN that has a higher level of error. than you would with a bug free implementatn.

$$gradApprox = (J(theta + \varepsilon) - J(theta - \varepsilon)) / (2\varepsilon)$$

for i=1 to n

    $\theta p = \theta$;

    $\theta p (i) = \theta p (i) + \varepsilon$;

    $\theta m = \theta$;

    $\theta m (i) = \theta m (i) - \varepsilon$;

    $gradApprox (i) = (J(\theta p) - J(\theta m)) / (2\varepsilon)$;

end

Check that

gradApprox $\approx$ DVec

   ↑ from numerical      └ from back.prop

1) Implement back prop to compute DVec (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$)

2) Implement numerical gradient check to compute grad Approx

3) Make sure they give similar values.

4) Turn off gradient check

5) Use back prop code for learning.

### 2.3. Random Initializatn:

→ Zero initializatn: $\theta_{ij}^{(l)} = 0$ for $\forall$ i, j, l.

for every training examples $a_1^{(2)} = a_2^{(2)}$

$$\frac{\partial}{\partial \theta_{01}^{(1)}} J(\theta) = \frac{\partial}{\partial \theta_{02}^{(1)}} J(\theta)$$

$$\delta_1^{(2)} = \delta_2^{(2)}$$

updated weights would be equal to each other

$$\downarrow$$

$$\longleftarrow \theta_{01}^{(1)} = \theta_{02}^{(1)}$$

⇒ After each update, paramtrs corresponding to inputs going into each of two hidden units are identical.   $a_1^{(2)} = a_2^{(2)}$

⇒ $a_1^{(2)} = a_2^{(2)} = a_3^{(2)} = \ldots = a_n^{(2)}$  means that ∀ of our hidden units are computing the exact same func. of the input! (highly redundant = equivalent with a unique feature, one unit on that layer).

INSTEAD!

→ Initialize each $\theta_{ij}^{(l)}$ to a random value in $[-\varepsilon, \varepsilon]$

  Theta1 = rand (10,11) * (2ε) − ε;    ⟶  $-\varepsilon < \text{Theta1} < \varepsilon$

  Theta 2 = rand (1,11) * (2ε) − ε;    ⟶  $-\varepsilon < \text{Theta2} < \varepsilon$

## 2.4. Putting It Together:

→ The more hidden units → the better.

Training a NN.
1. Randomly initialize weights
2. FP  ($h_\theta(x^{(i)})$)
3. Compute $J(\theta)$
4. BP  ($\frac{\partial}{\partial\theta_{jk}^{(l)}} J(\theta)$)

⎰ for i=1 to m
⎱     FP & BP  using $(x^{(i)}, y^{(i)})$
        (Get activatns $a^{(l)}$ & deltas $\delta^{(l)}$ for $l = 2,3,\ldots,L$)

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

     compute $\frac{\partial}{\partial\theta_{jk}^{(l)}} J(\theta)$

5. use grad. checking.
   then disable checking code.
6. Use grad. desc. or advanced optimizatn method with backprop. to try to minimize $J(\theta)$ as a fnn of θ.

    $J(\theta)$ — "non-convex", global minimum is not guaranteed, but in practice this is not a huge problem.

❗ what will if we try to maximize $J(\theta)$ at first, then initialize the weights with those (maximized $J(\theta)$), in order to, hopefully, end up with global minimum, Or at least, with better local minimum ?❗

→ Automous driving.

Once every 2 seconds, ALVINN digitizes a video img of the road ahead, and records the person's steering drn.

(3 layered Network)

→This same procedure is repeated for other road types