# Neural and brain

Neural network . thuật toán bắt chước bộ não
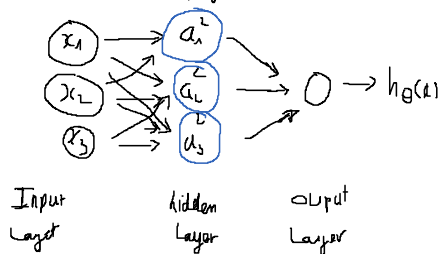


Input          hidden          Ouput

neuron model: logistic unit



$$h_{\theta(x)} = \frac{1}{1+e^{-\theta^T x}}$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_1 \\ x_n \end{bmatrix} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_n \end{bmatrix}$$

weight
(parameter)

Sigmoid (Logistic) activation function

$$a(z) = \frac{1}{1+e^{-z}} \quad , z = \theta^T x$$

Neural network



Input       hidden       ouput
Layer       Layer        Layer

$a_i^j$ : activation   unit $i$ of layer $j$

$\theta$ matrận trong số kiểm soát hàm
ảnh xạ từ một layer

$$a_1^{(2)} = g\left(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3\right)$$

$$a_2^{(2)} = g\left(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3\right)$$

$$a_3^{(2)} = g\left(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3\right)$$

$$h_\theta(x) = a_1^{(3)} = g\left(\theta_{10} a_0^{(2)} + \theta_{11} a_1^{(2)} + \theta_{12} a_2^{(2)} + \theta_{13} a_3^{(2)}\right)$$

nếu mạng có $S_j$ unit ở layer $j$, $S_{j+1}$ ở layer $j+1$, thì
ma trận $\theta^{(j)}$ có kích cỡ là $S_{j+1} \times (S_j + 1)$

---

Vectơ hóa công thức

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \qquad Z^{(j)} = \begin{bmatrix} z_1^{(j)} \\ z_2^{(j)} \\ z_3^{(j)} \end{bmatrix}$$

$$Z^{(2)} = \theta^{(1)} x$$

$$a^2 = g(Z^2)$$

$$R^3 \qquad R^3$$

for ward Propagation

$R^3$ $R^3$

thém $a_0^2 = 1$

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

$h_\theta(x) = a^3 = g(z^3)$

Forward Propagation

$x_1$ and $x_2$

true table

| $x_1$ | $x_2$ | $x_1$ and $x_2$ | |
|---|---|---|---|
| 1 | 1 | 1 | $y(30)=0$ |
| 1 | 0 | 0 | $y(-10)=0$ |
| 0 | 1 | 0 | $y(-10)\approx 0$ |
| 0 | 0 | 0 | $y(10)\approx 1$ |



$-30$, $20$, $20$ — $h_\theta(x)$

(not $x_1$) and (not $x_2$)

| $x_1$ | $x_2$ | $p$ |
|---|---|---|
| 1 | 1 | 0 |
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 0 | 1 |



$+10$, $-20$, $-20$ — $h_\theta(1)$

$x_1$ or $x_2$

| $x_1$ | $x_2$ | $p$ |
|---|---|---|
| 1 | 1 | 1 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 0 | 0 | 0 |



$-10$, $20$, $10$ — $h_\theta(x)$



$-30$, $20$, $10$, $-20$, $+10$, $-20$

$-10$, $20$, $20$ — $h_\theta(x)$

$x_1$ X NOR $x_2$

| $x_1$ | $x_2$ | $p$ |
|---|---|---|
| 1 | 1 | 1 |
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 0 | 1 |

Multiclass

multiple    Output    One vs all



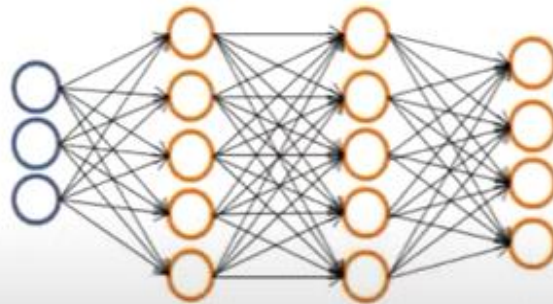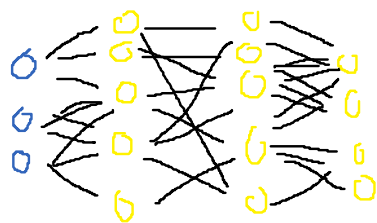Pedestrian        Car        Motorcycle        Truck

$h_\Theta(x) \in \mathbb{R}^4$

$h_\Theta(x) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$   Pedestrian

. . .

$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$   Car

training set  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \ldots (x^{m}, y^{m})$

$y$ one of  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

Cost function



$\{(x^1, y^1), (x^2, y^2), (x^3, y^3), \ldots, (x^m, y^m)\}$ m example

$L = $ no of Layer

$S_\ell = $ no of unit not Include bias

In layer $\ell$

binary classification

$y = 0$ or $1$

multi-class classification K class

$y \in \mathbb{R}^K$. $\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \ldots \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$

1 ouput unit

$y \in \mathbb{R}$

$S_L = 1$

K output units

$S_L = K$

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

$$+ \frac{\lambda}{2m} \sum_{\ell=1}^{\ell-1} \sum_{i=1}^{S_\ell} \sum_{j=1}^{S_{\ell+1}} (\theta_{ij}^{(\ell)})^2$$

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

$(h_\theta(x))_i = i^{th}$ ouput

Cost function

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

$$+ \frac{2\lambda}{m}\sum_{l=1}^{L-1}\sum_{i}^{S_l}\sum_{j}^{S_{l+1}}\left(\theta_{ij}^L\right)^2$$

min $J(\theta)$

$\theta$

$\Delta_j^l =$ error of node $j$ in layer $l$

$L = 4$

$\Delta_j^4 = a_j^4 - y_j$

$\Delta_j^3 = (\theta^3)^T \Delta_j^4 *. g'(z^3)$    <span style="color:blue">$a^3 *(1-a^3)$</span>   <span style="color:red">VD $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} *. \begin{bmatrix} 2 & 2 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 4 \\ 3 & 4 \end{bmatrix}$</span>

$\Delta_j^2 = (\theta^2)^T \Delta_j^3 *. g'(z^2)$    <span style="color:blue">$a^2 *(1-a^2)$</span>

<span style="color:red">note: $*.$ là nhân từng phần tử của ma trận</span>

<span style="color:blue">$\frac{\partial}{\partial\theta_{ij}}J(\theta) = a_j^l \cdot \Delta_I^{l+1}$</span>

Algorith

training set $(x^1, y^1), (x^2, y^2), (x^3, y^3), \dots, (x^4, y^4)$

Set $\Delta_{ij}^l = 0$ for all $i, j, l$

For $i = 1$ to $m$

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

    Perform forward propagation to compute
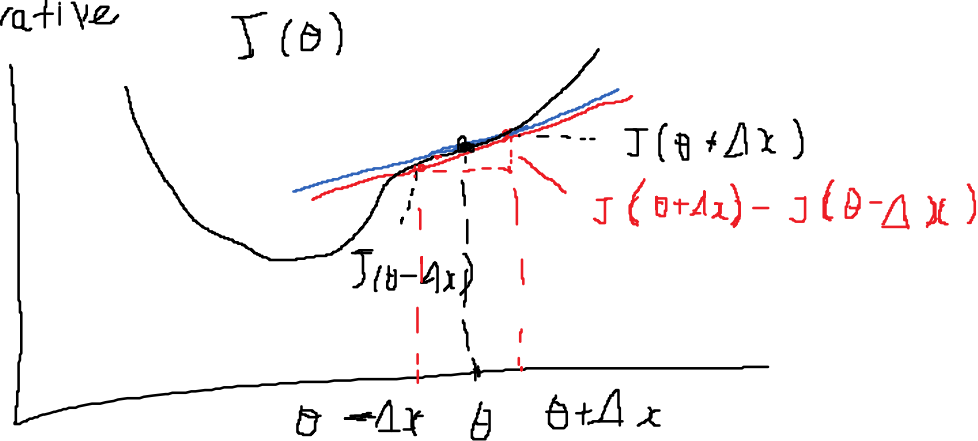        $a^l$ for $l = 1, 2, 3, \dots, L$

    Compute $\Delta^L = a^L - y^i \dots$

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

$\rightarrow \quad \Delta_{ij}^{\ell} = \Delta_{ij}^{\ell} + a_j^{\ell} \Delta_i^{\ell+1}$

$\underline{\Delta}^{\ell} = \Delta^{\ell} + \Delta^{(\ell+1)} \cdot a^{(\ell)T}$

derivative

$J(\theta)$



$J(\theta + \Delta x)$

$J(\theta + \Delta x) - J(\theta - \Delta x)$

$J(\theta - \Delta x)$

$\theta - \Delta x \quad \theta \quad \theta + \Delta x$

$$\frac{\partial J(\theta)}{\partial \theta} = \frac{J(\theta + \Delta x) - J(\theta - \Delta x)}{2 \Delta x}$$

$\Delta x = 10^{-4} = 0$

Voi $\quad J(\theta); \theta = \left[ \theta_1, \theta_2, \theta_3, \theta_4, \ldots, \theta_n \right]$

$$\frac{\partial J(\theta)}{\partial \theta_1} = \frac{J(\theta_1 + \Delta x, \theta_2, \theta_3, \ldots, \theta_n) - J(\theta_1 - \Delta x, \theta_2, \ldots, \theta_n)}{2 \Delta x}$$

$$\frac{\partial J \theta}{\partial \theta_2} = \frac{J(\theta_1, \theta_2 + \Delta x, \theta_3, \theta_4, \ldots, \theta_n) - J(\theta_1, \theta_2 - \Delta x, \ldots, \theta_n)}{2 \Delta x}$$

$\ldots$

$$\frac{\partial J(\theta)}{\partial \theta_n} = \frac{J(\theta_1, \theta_2, \ldots, \theta_n + \Delta x) - J(\theta_1, \theta_2, \ldots, \theta_n - \Delta x)}{2 \Delta x}$$

Octave

```
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;
```

Step

use  back propagation  to  compute  Dvec

use  Gradient checking  to  compute  Gradapprox

Make sure   Dvec ≃ Gradapprox

Disable  Gradient checking  then  run  back propagation

# Random initialization

If the dimensions of Theta1 is 10x11, Theta2 is 10x11 and Theta3 is 1x11.

matran 10x11 random(0,1)

```
Theta1 = rand(10,11) * (2 * INIT_EPSILON) - INIT_EPSILON;
Theta2 = rand(10,11) * (2 * INIT_EPSILON) - INIT_EPSILON;
Theta3 = rand(1,11) * (2 * INIT_EPSILON) - INIT_EPSILON;
```

Cách để chọn cấu trúc net work

No of input : Dimension of feature $x^{(1)}$

No of output : Number of classes

hidden layer : default 1, or $>1$ layer, have same no of hidden unit (càng nhiều càng tốt)

Các bước training

1 random Intialization các trọng số $\theta$ (weigt)

2. dùng forward propagation để tính $h_\theta(x^{(i)})$ cho $y^{(i)}$

3. xây dựng code tính cost function $J(\theta)$

4. khai triển back prop để tính đạo hàm $\dfrac{\partial J(\theta)}{\partial \theta_{ij}}$

5 dùng Gradient checking để kiểm tra

6. dùng Gradient descent để minimize $J(\theta)$