# Logistic Regression



$$z = \theta^T x \quad \Big| \quad P(y=1 \mid x, \theta) = h_\theta(x) = \frac{1}{1+e^{-z}}$$

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

$$\hat{y} = \begin{cases} 1 & , \text{if } h_\theta(x) \geq 0.5 \\ \\ 0 & , \text{if } h_\theta(x) < 0.5 \end{cases}$$

$$\hat{y} = \begin{cases} 1 & , \text{if } z \geq 0 \\ \\ 0 & , \text{if } z < 0 \end{cases}$$

**\[\*\*\] cost function for $i^{th}$ sample: (Binary cross entropy loss)**

$$J_i(\theta) = \begin{cases} -\log_e(h_\theta(x)) & , \text{if } y=1 \\ \\ -\log_e(1-h_\theta(x)) & , \text{if } y=0 \end{cases}$$

$$\hookrightarrow J(\theta) = \frac{1}{n} \sum_{i=1}^{n} \left[ y^i \times \log_e(h_\theta(x^i)) + (1-y^i)\log_e(1-h_\theta(x^i)) \right]$$
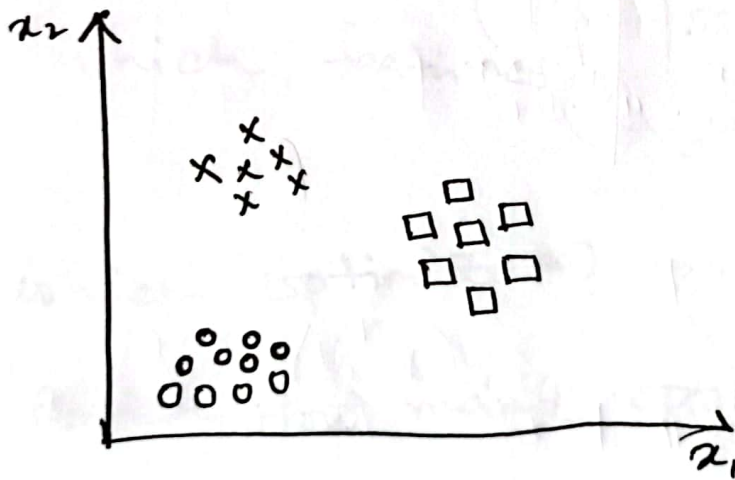
$$\underline{P1}$$

$$\underline{Rev.}$$

$$P2$$

## ✺ Multiclass classification

$$
\overset{0}{Email} : \overset{0}{Friends} / \overset{1}{Colisue} / \overset{2}{spam}
$$

$$
\underset{0}{image:} \underset{0}{human} / \underset{1}{horse} / \underset{2}{cat} / \underset{3}{dog}
$$



$$y = \{0, 1, 2\}$$

$$o \quad x \quad \square$$

LR এর extention

$y = \{0, 1, 2\}$

**\*\*** <u>Algo-1</u> | <u>One vs all</u>

<u>idea:</u> Train a logistic regression classifier $h_{\theta_i}^{(i)}$ for each class $i$ to predict the probability of $\boxed{y = i}$

3 টে $\theta$ create করবো।

$(h_{\theta_0}^{(0)}) = \theta_1^T x$  $(h_q^{(1)}) = \theta_1^T x$  $(h_{\theta_2}^{(2)}) = \theta_2^T x$

$\{\theta_{00}, \theta_{01}, \ldots, \theta_{0d}\}, \{\theta_{10}, \theta_{11}, \ldots, \theta_{1d}\}, \{\theta_{20}, \theta_{21}, \ldots, \theta_{2d}\}$

এটা input ফিলে <u>output দিবে যে</u> ~~এটার data x~~ এই টা data x '0' হওয়ার <u>probability</u> <u>কত</u>

এর input ফিলে <u>যেকোনো Data x '1' হওয়ার</u> <u>probability</u> <u>কত</u>

<u>Data x এর</u> <u>'2' হওয়ার</u> <u>probability</u> <u>কত</u>

so $\left[ h_{\theta_0}^{(0)} \cdots = 0.2 \right]$

$h_{\theta_1}^{(1)} = 0.1$

$h_{\theta_2}^{(2)} = 0.9$

যেটার probability বেশি
Data x @ class এ
<u>হবে।</u>

$$y_{pred} = \text{argmax}\; h_{\theta_i}^{(i)}(x)$$

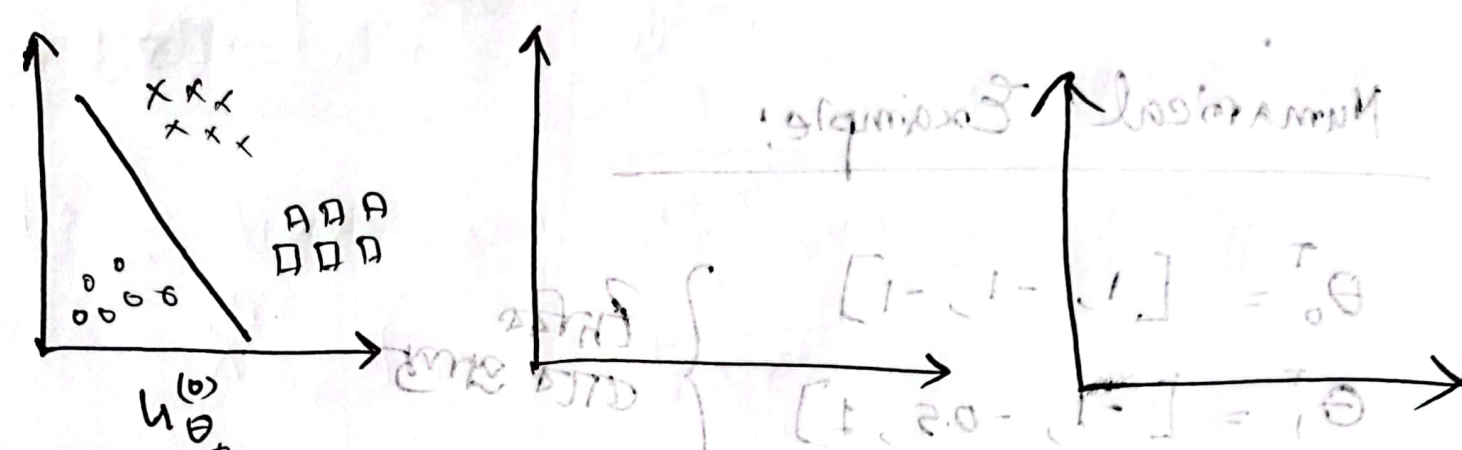$$h_{\theta_0}^{(0)}(x) = g(\theta_0^T x) = \frac{1}{1+e^{-\theta_0^T x}} = P(y=0 \mid x;\theta_0)$$

$$[\theta_{00}, \theta_{01}, \theta_{02}, \ldots, \theta_{0d}]$$

$$h_{\theta_1}^{(1)}(x) = g(\theta_1^T x) = \frac{1}{1+e^{-\theta_1^T x}} = P(y=1 \mid x;\theta_1)$$

$$[\theta_{10}, \theta_{11}, \theta_{12}, \ldots, \theta_{1d}]$$

$$h_{\theta_2}^{(2)}(x) = g(\theta_2^T x) = \frac{1}{1+e^{\theta_2^T x}} = P(y=2 \mid x;\theta_2)$$
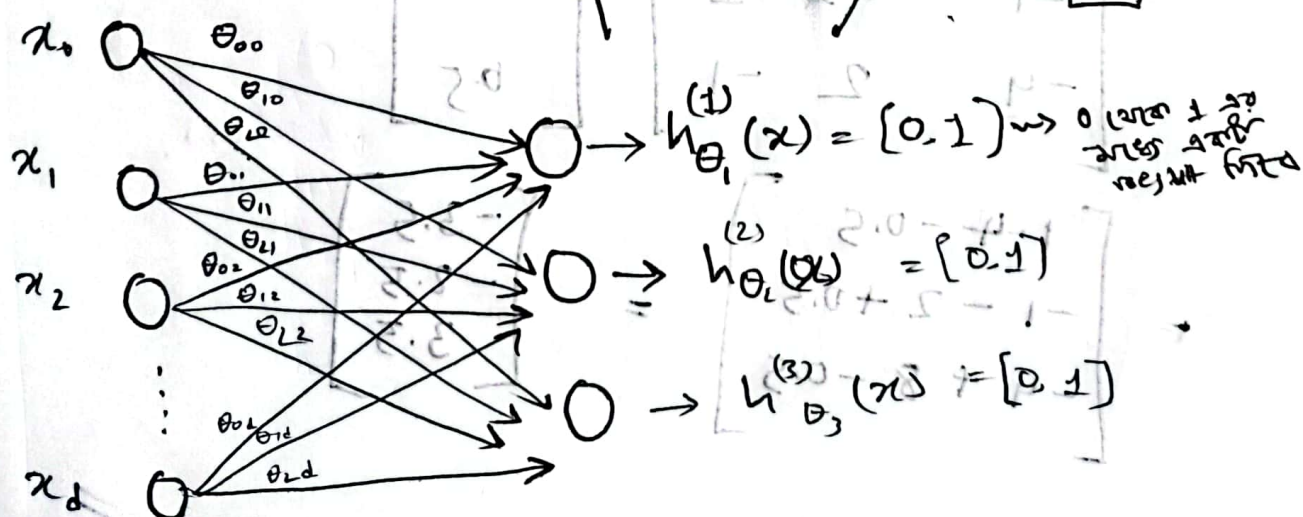
$$[\theta_{20}, \theta_{21}, \theta_{22}, \ldots, \theta_{2d}]$$

vectorized:

$$W = \begin{bmatrix} \Theta_0^T \\ \Theta_1^T \\ \Theta_2^T \end{bmatrix} = \begin{bmatrix} \Theta_{00}, \Theta_{01}, \cdots, \Theta_{0d} \\ \Theta_{10}, \Theta_{11}, \cdots, \Theta_{1d} \\ \Theta_{20}, \Theta_{22}, \cdots, \Theta_{2d} \end{bmatrix}$$

$$h_\theta(x) = g(Wx) = g\left(\begin{bmatrix} \Theta_0^T x \\ \Theta_1^T x \\ \Theta_2^T x \end{bmatrix}\right) = \begin{bmatrix} g(\Theta_0^T x) \\ g(\Theta_1^T x) \\ g(\Theta_2^T x) \end{bmatrix}$$



$$h_{\theta_1}^{(1)}(x) = [0.1] \rightsquigarrow$$ 0 ज्यादा 1 कम

$$h_{\theta_2}^{(2)}(x) = [0.1]$$

$$h_{\theta_3}^{(3)}(x) = [0.1]$$

## Numarical Example:

$$\theta_0^T = [1, -1, -1]$$

$$\theta_1^T = [-1, -0.5, 1]$$

$$\theta_2^T = [-4, 2, -1]$$

$$x = \begin{bmatrix} 1 \\ 4 \\ 0.5 \end{bmatrix}$$

$$W = \begin{bmatrix} 1 & -1 & -1 \\ -1 & -0.5 & 1 \\ -4 & 2 & -1 \end{bmatrix}$$

$$\therefore \ z = Wx$$

$$= \begin{bmatrix} 1 & -1 & -1 \\ -1 & -0.5 & 1 \\ -4 & 2 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 4 \\ 0.5 \end{bmatrix}$$

$$= \begin{bmatrix} 1 - 4 - 0.5 \\ -1 - 2 + 0.5 \\ -4 + 8 - 0.5 \end{bmatrix} = \begin{bmatrix} -3.5 \\ -2.5 \\ 3.5 \end{bmatrix}$$

**[\*\*] Problems - of one vs all classifier:**

① Have to train seperately

② No probabilistic interpretation for $h_\theta$ [sum($h$) > 1]

$$\left( \begin{array}{l} \text{অবশ্যই sum = 1 হয় না,} \\ \text{যদি each class উর probability নিয়ে} \\ \text{তাহলে sum = 1 হবে} \end{array} \right)$$

**[\*\*] Solution:** $\begin{cases} ① \text{ we want sum($h$) = 1} \\ ② \text{ we want end to end learning} \end{cases}$

3টা seperate classifier train করবো না, একটা বড়ো একটা single অবশ্যই classifier বানাই!!

<mark>Softmax classifier</mark>

⇒ **modification to the Previous algorithm:**

① Sigmoid function to softmax function

② change $y$ to one-hot vector.

③ from $J(\theta) = \underset{\text{Cross entropy}}{\underline{BCE}}$ to $\widetilde{CE}$

Binary → BCE

Cross Entropy → CE

① $\Rightarrow$

$$Wx = [z_0, z_1, \ldots z_{K-1}]$$

$\underset{(K \times d+1)}{\downarrow} \underset{((d+1) \times 1)}{\downarrow}$

$\rightarrow (K \times 1) K$

$$\text{sigmoid} (Wx)^T = \left[ \frac{1}{1+e^{-z_0}}, \frac{1}{1+e^{-z_1}}, \ldots, \right]$$

$\rightarrow$ Softmax function, $f_i = \dfrac{e^{z_i}}{\sum\limits_{i=0}^{K-1} e^{z_i}}$

$$\text{Softmax} (Wx)^T = [f_0, f_1, \ldots, f_{K-1}]$$

Probability $\begin{cases} \left( f_1 = \dfrac{e^{z_1}}{e^{z_0} + e^{z_1} + e^{z_2}} \right) \\[4mm] \left( \sum\limits_{i=0}^{K-1} f_i = 1 \right) \end{cases}$

②  ⟹  one hot ~~vector~~ encoding

$$y = \{0, 1, 2\} = \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\}$$ ①

③  ⟹

Cross entropy loss,

$$J(w) = -\frac{1}{n} \sum_{i=1}^{n} \underbrace{\sum_{k=0}^{k-1} y_k^{(i)} \log f_k}_{\;}$$

$\underbrace{\phantom{xxxxxxxx}}$ avg over all training samples

$J_i(w)$

Optimizer ⟹ G.P

$w^0 = $ random

while not converged:

$$w^{(k+1)} = w^k - \nabla_w J$$

if $\nabla_w J \approx 0$:

converged == True

$$Wx$$



$$g(z) \text{ [Softmax]}$$

$$\begin{bmatrix} \theta_{00} & \theta_{01} & \cdots & \theta_{0d} \\ \theta_{10} & \theta_{11} & \cdots & \theta_{1d} \\ \theta_{20} & \theta_{21} & \cdots & \theta_{2d} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} = \begin{bmatrix} z_0 \\ z_1 \\ z_2 \end{bmatrix} \xrightarrow{\text{softmax}} \begin{bmatrix} f_0 \\ f_1 \\ f_2 \end{bmatrix}$$

$$\theta_0 \qquad \theta_1 \qquad \theta_d$$

$$(3,d) \times (d,1) \qquad (3\times1)$$

$$\theta_0 = [\,0,\quad 0.01,\quad -0.05,\quad 0.1,\quad 0.05\,]$$

$$\theta_1 = [\,0.2,\quad 0.7,\quad 0.2,\quad 0.05,\quad 0.16\,]$$

$$\theta_2 = [\,-0.3,\quad 0.0,\; -0.45,\; -0.2,\quad 0.33\,]$$

$$x_0 = [\,1,\quad 15,\quad 22,\quad 44,\quad 56\,], \quad y = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

(Matrix form में आएगा) $z = [Wx]$

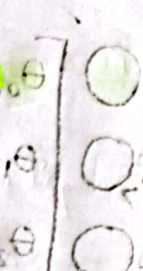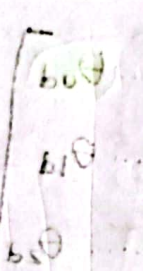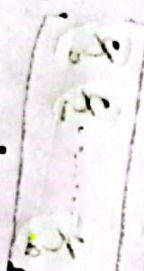$$z_0 = 0\times1 - 0.01\times15 - 0.05\times22 - 0.1\times44 + 0.05\times56$$
$$= -2.58$$

$$z_1 = 0.02\times1 - 0.7\times15 + 0.2\times22 - 0.05\times44 + 0.16\times56$$
$$= 0.86$$

$$z_2 = -0.3\times1 - 0\times15 - 0.45\times22 + 0.2\times44 + 0.33\times56$$
$$= 0.28$$

$$\therefore \begin{bmatrix} z_0 \\ z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} -2.58 \\ 0.86 \\ 0.28 \end{bmatrix}$$

$$\therefore f\left(\begin{bmatrix} z_0 \\ z_1 \\ z_2 \end{bmatrix}\right) = \begin{bmatrix} \dfrac{e^{-2.58}}{e^{-2.58}+e^{0.86}+e^{0.28}} \\[2mm] \dfrac{e^{0.86}}{e^{-2.58}+e^{0.86}+e^{0.28}} \\[2mm] \dfrac{e^{0.28}}{e^{-2.58}+e^{0.86}+e^{0.28}} \end{bmatrix} = \begin{bmatrix} 0.016 \\ 0.631 \\ 0.353 \end{bmatrix} \approx \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$y_{pred}$$

Now, $y = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$ , $y_{pred} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \approx \begin{bmatrix} 0.016 \\ 0.631 \\ 0.323 \end{bmatrix}$

$$\therefore\ J_d^{(w)} = -\sum_{k=0}^{K-1} y_k \log f_k$$

$$= -[y_0 \times \log f_0 + y_1 \times \log f_1 + y_2 \times \log f_2]$$

$x_0$ এর জন্য error

$$= -(0 \times \log(0.016) + 0 \times \log(0.631) + 1 \times \log(0.323))$$

$$= -\log(0.323) = 1.04$$

(Suppose)

$\therefore\ x_1$ এর জন্য, $J_1(w) = 2.06$

$\therefore\ x_2$ এর জন্য, $J_2(w) = 0.08$

$\therefore\ x_3$ এর জন্য, $J_3(w) = 1.34$

$$\therefore\ J(w) = \frac{1}{n}\sum_{i=0}^{n-1}\left(-\sum_{k=0}^{K-1} y_k \log f_k\right)$$

$n \rightarrow$ number of samples

$K \rightarrow$ number of class

$$= \frac{1}{4}[J_0(w) + J_1(w) + J_2(w) + J_3(w)]$$

$$= \frac{1}{4}(1.04 + 2.06 + 0.08 + 1.34) = 1.13$$

# Practical consideration

① (Numbers of neurron
  ⤷ over fitting
  ⤷ need to tune
   the hyper parametene
    ⤷ need validation set

② over fiting sol$^n$ : Regularization

$$J(\theta) + \lambda \sum |\theta|^2$$

③ activation function:
  Sigmoid   saturates $\} \rightarrow$ ( y is) +ve
  ⤷ use $\quad$ tanh

     $\rightarrow z$, if $z \geqslant 0$
** Relu(z) $\{$
     $\rightarrow 0$, other wise