# SOFTMAX REGRESSION

## I. MOTIVATION



| width | Label |
|-------|-------|
| 1.4 | 0 |
| 1.5 | 0 |
| 4.5 | 1 |
| 5.6 | 1 |

$x$ width $\to$ $O$ $\xrightarrow{w_0}$ $(z_1)$ $\to$ sigmoid $\to$ $p$

$1$ $\xrightarrow{b_0}$

1 & 0 riêng thì ?

$\to$ Tại sao phân loại 2 cái mà chỉ có 1 node ? Nếu cho 2 node đại diện

$x$ width $\to$ $O$ $\xrightarrow{w_0}$ $w_1$ $\xrightarrow{b_0}$ $b_1$ $(z_0)$ $(z_1)$ $\to$ ? $<$ $P(y=0|x)$ / $P(y=1|x)$

softmax

- Softmax function : $P_i = f(z_i) = \dfrac{e^{z_i}}{\sum_j e^{z_j}}$

Tại sao không dùng xs bt $\dfrac{z_i}{\sum z_i}$

để tránh tuyến tính $\to$ k nhạy cảm

$\to$ ổn định hơn với những số cực nhỏ hoặc cực lớn

+ Note: khi cài đặt softmax trên máy dễ bị tràn $\to$ trừ đi max để scale lại

- Loss function : $-y \log \hat{y}_1 - (1-y)\log \hat{y}_0$

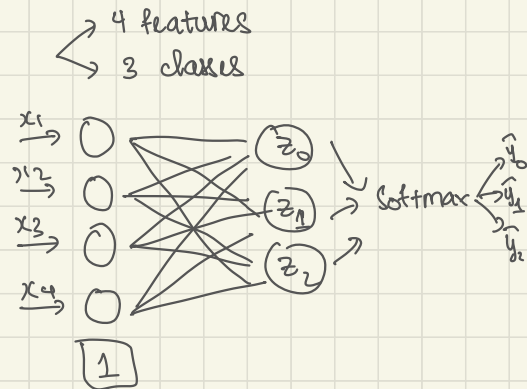## II. MODEL CONSTRUCTION

Feature Label

| width | Label |
|-------|-------|
| 1.4 | 0 |
| 1.5 | 0 |
| 4.5 | 1 |
| 5.6 | 1 |

class : 2 (0 & 1) $\to$ số lượng output node

feature : 1 $\to$ số lượng input node

| S_length | S_width | P_length | P_width | label |
|---|---|---|---|---|
| — | — | — | | 0 |
| — | — | — | | 1 |
| — | — | | | 0 |
| — | — | | | 1 |
| — | | | — | 2 |
| — | | | | 2 |

4 features
3 classes

$x_1$ ○
$x_2$ ○     $z_0$
$x_3$ ○     $z_1$ → Softmax → $\hat{y}_0$, $\hat{y}_1$, $\hat{y}_2$
$x_4$ ○     $z_2$
1

$3 \times 5 = 15$ parameters

VD :

folder
├── $t_1$ : cat
├── $t_2$ : dog
├── $t_3$ : duck
├── $t_4$ : tiger
└── $t_5$ : goat

image : $10 \times 10$

feature : 100
class : 5

# III. LOSS FUNCTION

$x$
petal length → ○ → $z_0$ → Softmax → $\hat{y}_0$
1 → $z_1$ → → $\hat{y}_1$

$$z_0 = x w_0 + b_0$$
$$z_1 = x w_1 + b_1$$

$$z = \theta^T x \quad \left( z = \begin{bmatrix} b_0 & w_0 \\ b_0 & w_1 \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix} = \begin{bmatrix} \theta_0^T \\ \theta_1^T \end{bmatrix} \begin{bmatrix} 1 \\ x \end{bmatrix} = \theta^T x \right)$$

$$y_0 = \frac{e^{z_0}}{\sum_j e^{z_j}}$$
$$y_1 = \frac{e^{z_1}}{\sum_j e^{z_j}}$$

$$y = \frac{e^z}{\sum_j e^{z_j}}$$

$$L(\theta) = -y_0 \log \hat{y}_0 - y_1 \log \hat{y}_1$$

- One-hot encoding:

code : $y = [\ 0\ ,\ 0\ ]$

$$\frac{\partial L}{\partial z_0} = \frac{\partial L}{\partial \hat{y}_0} \cdot \frac{\partial y_0}{\partial z_0} + \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_0}$$

$$\frac{\partial L}{\partial z_1} =$$

$y[y] = 1 \quad \xrightarrow{y=0} \quad y = [\ 1\ ,\ 0\ ]$

$y = 1$

$y = [\ 0\quad 1\ ]$
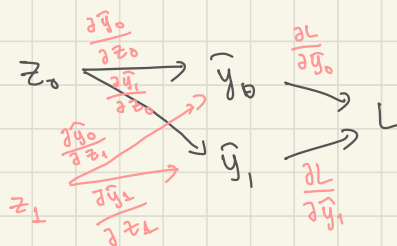
$$L(\theta) = -y \log \hat{y}_0 - y_1 \log \hat{y}_1 = \sum_{i=0}^{1} y_i \log \hat{y}_i = -y^T \log \hat{y}$$

$$\left(\frac{u}{v}\right)' = \frac{u'v - v'u}{v^2}$$

$$y_0 = \frac{e^{z_0}}{\sum_{j=0}^{1} e^{z_j}} \qquad y_1 = \frac{e^{z_1}}{\sum_{i=0}^{1} e^{z_i}}$$

$$\left(\frac{1}{u}\right)' = -\frac{1}{u^2} \cdot u'$$

$$(e^u)' = u' \cdot e^u$$

$$\frac{\partial L}{\partial \hat{y}_0} = \frac{-y_0}{\hat{y}_0} \qquad \frac{\partial L}{\partial \hat{y}_1} = \frac{-y_1}{\hat{y}_1}$$

$z_0 \ \xrightarrow{\frac{\partial \hat{y}_0}{\partial z_0}} \ \hat{y}_0 \ \xrightarrow{\frac{\partial L}{\partial \hat{y}_0}} \ L$

$z_1 \ \xrightarrow{\frac{\partial \hat{y}_1}{\partial z_1}} \ \hat{y}_1 \ \xrightarrow{\frac{\partial L}{\partial \hat{y}_1}}$

$$\frac{\partial \hat{y}_0}{\partial z_1} = \frac{-e^{z_0}}{(e^{z_0}+e^{z_1})^2} \cdot e^{z_1} = \frac{-e^{z_0}}{(e^{z_0}+e^{z_1})} \cdot \frac{e^{z_1}}{(e^{z_0}+e^{z_1})}$$

$$= -\hat{y}_0 \cdot \hat{y}_1$$

$$\frac{\partial \hat{y}_0}{\partial z_0} = \left(\frac{e^{z_0}}{e^{z_0}+e^{z_1}}\right)' = \frac{e^{z_0} \cdot \sum - e^{z_0} \cdot e^{z_0}}{(e^{z_0}+e^{z_1})^2} = \frac{e^{z_0}(\sum - e^{z_0})}{(\sum^2)} = \frac{e^{z_0}}{\sum} \cdot \left(\frac{\sum}{\sum} - \frac{e^{z_0}}{\sum}\right)$$

$$= \hat{y}_0 \cdot (1 - \hat{y}_0)$$

$$\frac{\partial y_1}{\partial z_1} = \hat{y}_1 (1 - \hat{y}_1)$$

$$\frac{\partial y_1}{\partial z_0} = -\hat{y}_0 \cdot \hat{y}_1$$

Derivative
$$\frac{\partial \hat{y}_i}{\partial z_j} = \begin{cases} y_i(1 - y_i) & i = j \\ -\hat{y}_i\, y_j & i \neq j \end{cases}$$

$$\frac{\partial L}{\partial z_0} = \frac{\partial L}{\partial \hat{y}_0} \cdot \frac{\partial y_0}{\partial z_0} + \frac{\partial L}{\partial \hat{y}_1} \cdot \frac{\partial \hat{y}_1}{\partial z_0} = -\frac{y_0}{\hat{y}_0} \cdot \hat{y}_0(1-\hat{y}_0) + \frac{y_1}{\hat{y}_1} \cdot \left( \hat{y}_0 \cdot \hat{y}_1 \right)$$

$$= -y_0(1-\hat{y}_0) + y_1 \cdot \hat{y}_0 = -y_0 + y\hat{y}_0 + y_1\hat{y}_0 = \hat{y}_0(y_0 + y_1) - y_0$$

$$\frac{\partial L}{\partial z_1} = \hat{y}_1 - y_1 \qquad\qquad = \hat{y}_0 - y_0 \quad \left(\begin{array}{l}\text{do one-hot} \\ \text{thì } y_1 = 0\end{array}\right)$$

$$\frac{\partial L}{\partial w_0} = x(\hat{y}_0 - y_0)$$

# IV. SUMMARY

**1. Forward computation**

$$z = \theta^T x$$

$$y = \frac{e^z}{\sum_{j=0}^{1} e^{z_j}}$$

**2. loss function**

$$L(\theta) = -y^T \log y$$

**3. Derivative**

$$\frac{\partial L}{\partial w_i} = x(\hat{y}_i - y) \qquad \frac{\partial L}{\partial b_i} = \hat{y}_i - y_i$$

$$\nabla_\theta L = x(\hat{y} - y)^T$$

**4. Update**

$$\theta = \theta - \eta L'_\theta$$

# V. A SUGGEST FUNCTION

data $\searrow$ # feature : 1
$\qquad\searrow$ # classes : 3

$$x \rightarrow \bigcirc \quad \bigcirc \rightarrow \quad \boxed{\text{softmax}} \quad \begin{array}{l} \rightarrow \hat{y}_0 \\ \rightarrow \hat{y}_1 \\ \rightarrow \hat{y}_2 \end{array}$$

$$L(\theta) = -\frac{y(1-y)}{-2}\log(y_2) - y(2-y)\log(\hat{y}_1) - (1-y)\left(\frac{2-y}{2}\right)\log(\hat{y}_0)$$

$$\underbrace{\qquad}_{y_2} \qquad \underbrace{\qquad}_{y_1} \qquad \underbrace{\qquad}_{y_0}$$

$\rightarrow$ Khá phức tạp nếu nhiều label $\rightarrow$ one-hot encoding

$$L(\Theta) = -y_2 \log(\hat{y}_2) - y_1 \log(\hat{y}_1) - y_0 \log(\hat{y}_0)$$

# VI. MODEL

$$\frac{\partial L}{\partial b_0} = \hat{y}_0 - y_0$$
$$\frac{\partial L}{\partial w_0} = x(\hat{y}_0 - y_0)$$

$x$

$$\frac{\partial L}{\partial b_1} = \hat{y}_1 - y_1$$
$$\frac{\partial L}{\partial w_1} = x(\hat{y}_1 - y_1)$$

$$\frac{\partial L}{\partial z_0} = \hat{y}_0 - y_0$$

$b_0$  $w_0$   $b_1$  $w_1$

$$\frac{\partial L}{\partial z_1} = \hat{y}_1 - y_1$$

$$z_0 = w_0 x + b_0$$

$$z_1 = w_1 x + b_1$$

$$\hat{y}_0 = \frac{e^{z_0}}{\sum_{i=0}^{1} e^{z_i}}$$

$$\hat{y}_1 = \frac{e^{z_1}}{\sum_{i=0}^{1} e^{z_i}}$$

$$L = -y_0 \log \hat{y}_0 - y_1 \log \hat{y}_1$$

$\leftarrow$  label  $\boxed{y}$  $\begin{bmatrix} one \\ hot \end{bmatrix}$