# Exercise 5 - Softmax regression

## Đặng Linh Anh

## January 2021

**1.** Xây dựng công thức tính forward và backward (tính đạo hàm cho từng tham số) cho bài toán softmax regression dùng phương pháp dựa trên hàm delta.

(a) Stochastic gradient descent

(b) Batch gradient descent

### Lời giải

Cho dataset có: $n$ features, $N$ samples, $k$ outputs

$$\boldsymbol{\Theta} = \begin{bmatrix} w_{11} & w_{21} & \cdots & w_{k1} \\ w_{12} & w_{22} & \cdots & w_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ w_{1n} & w_{2n} & \cdots & w_{kn} \\ b_1 & b_2 & \cdots & b_k \end{bmatrix} \tag{1}$$

(a) Stochastic gradient descent:
Cài đặt:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \\ 1 \end{bmatrix} \tag{2}$$

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_k \end{bmatrix} \tag{3}$$

$$\Delta = \begin{bmatrix} \delta_1 \\ \delta_2 \\ \vdots \\ \delta_k \end{bmatrix} \tag{4}$$

$$\delta_i = \begin{cases} 1, & \text{nếu } y = i \\ 0, & \text{khác} \end{cases} \tag{5}$$

i. Tính toán forward

$$\mathbf{z} = \boldsymbol{\Theta}^T \mathbf{x} \tag{6}$$

$$\hat{\mathbf{y}} = \frac{e^{\mathbf{z}}}{\sum e^{\mathbf{z}}} \tag{7}$$

1

Loss function:

$$L \ = \ \sum_{i=1}^{k} \delta_i \log \hat{y} \tag{8}$$

$$\Rightarrow L \ = \ \sum (\Delta \odot \log \hat{\mathbf{y}}) \tag{9}$$

ii. Tính toán backward

- Gradient cho từng tham số:

$$\frac{\partial L}{\partial \Theta_{ij}} \ = \ \sum_{u=1}^{k} \frac{\partial L}{\partial \hat{y}_u} \times \sum_{v=1}^{k} \frac{\partial \hat{y}_v}{\partial z_i} \times \frac{\partial z_i}{\partial \Theta_{ij}} \tag{10}$$

Trong đó:

$$\frac{\partial z_i}{\partial \Theta_{ij}} \ = \ x_j, \tag{11}$$

$$\sum_{v=1}^{k} \frac{\partial \hat{y}_v}{\partial z_i} \ = \ \sum_{v=1}^{k} \frac{\partial}{\partial z_i} \left( \frac{e^{z_v}}{\sum_{j=1}^{k} e^{z_j}} \right) \tag{12}$$

$$\Leftrightarrow \sum_{v=1}^{k} \frac{\partial \hat{y}_v}{\partial z_i} \ = \ \frac{e^{z_i}}{\sum_{j=1}^{k} e^{z_j}} - \sum_{v=1}^{k} \frac{e^{z_i} e^{z_v}}{\left( \sum_{j=1}^{k} e^{z_j} \right)^2} \tag{13}$$

$$\Leftrightarrow \sum_{v=1}^{k} \frac{\partial \hat{y}_v}{\partial z_i} \ = \ \sum_{v=1}^{k} \hat{y}_v (\delta_v - \hat{y}_i), \tag{14}$$

$$\sum_{u=1}^{k} \frac{\partial L}{\partial \hat{y}_u} \ = \ -\sum_{u=1}^{k} \frac{\delta_u}{\hat{y}_u} \tag{15}$$

Từ đó, ta tìm được gradient:

$$\frac{\partial L}{\partial \Theta_{ij}} \ = \ -\sum_{u=1}^{k} \frac{\delta_u}{\hat{y}_u} \times \sum_{v=1}^{k} \hat{y}_v (\delta_i - \hat{y}_i) \times x_j \tag{16}$$

$$\Leftrightarrow \frac{\partial L}{\partial \Theta_{ij}} \ = \ x_j (\hat{y}_i - \delta_i) \tag{17}$$

Vectorization:

$$\frac{\partial L}{\partial \Theta} = \mathbf{x}^T (\hat{\mathbf{y}} - \Delta) \tag{18}$$

- Cập nhật tham số:

$$\mathbf{\Theta} := \mathbf{\Theta} - \eta \times \mathbf{x}^T (\hat{\mathbf{y}} - \Delta) \tag{19}$$

(b) Batch gradient descent

Cài đặt:

$$\mathbf{X} = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \cdots & x_1^{(m)} \\ x_2^{(1)} & x_2^{(2)} & \cdots & x_2^{(m)} \\ \vdots & \vdots & \ddots & \vdots \\ x_n^{(1)} & x_n^{(2)} & \cdots & x_n^{(m)} \\ 1 & 1 & \cdots & 1 \end{bmatrix} \tag{20}$$

$$\mathbf{Y} = \begin{bmatrix} y^{(1)} & y^{(2)} & \cdots & y^{(m)} \end{bmatrix} \tag{21}$$

$$\hat{\mathbf{Y}} = \begin{bmatrix} \hat{y}_1^{(1)} & \hat{y}_1^{(2)} & \cdots & \hat{y}_1^{(m)} \\ \hat{y}_2^{(1)} & \hat{y}_2^{(2)} & \cdots & \hat{y}_2^{(m)} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{y}_k^{(1)} & \hat{y}_k^{(2)} & \cdots & \hat{y}_k^{(m)} \end{bmatrix} \tag{22}$$

$$\Delta = \begin{bmatrix} \delta_1^{(1)} & \delta_1^{(2)} & \cdots & \delta_1^{(m)} \\ \delta_2^{(1)} & \delta_2^{(2)} & \cdots & \delta_2^{(m)} \\ \vdots & \vdots & \ddots & \vdots \\ \delta_k^{(1)} & \delta_k^{(2)} & \cdots & \delta_k^{(m)} \end{bmatrix} \tag{23}$$

$$\delta_i^{(u)} = \begin{cases} 1, & \text{nếu } y^{(u)} = i \\ 0, & \text{khác} \end{cases} \tag{24}$$

i. Tính toán forward

$$\mathbf{Z} = \mathbf{\Theta}^T \mathbf{x} \tag{25}$$

$$\hat{y}^{(u)} = \frac{e^{\mathbf{z}^{(u)}}}{\sum e^{\mathbf{z}^{(u)}}} \tag{26}$$

$$\Rightarrow \hat{\mathbf{Y}} = \begin{bmatrix} \hat{\mathbf{y}}^{(1)} & \hat{\mathbf{y}}^{(2)} & \cdots & \hat{\mathbf{y}}^{(u)} \cdots \hat{\mathbf{y}}^{(m)} \end{bmatrix} \tag{27}$$

Loss function:

$$L = \frac{1}{m} \sum_{u=1}^{m} \sum_{i=1}^{k} \delta_i^{(u)} \log \hat{y}_i^{(u)} \tag{28}$$

$$\Rightarrow L = \frac{1}{m} \underbrace{\begin{bmatrix} 1 & \cdots & 1 \end{bmatrix}}_{k} (\Delta \odot \log \hat{\mathbf{y}}) \underbrace{\begin{bmatrix} 1 & \cdots & 1 \end{bmatrix}^T}_{m} \tag{29}$$

ii. Tính toán backward
   - Gradient cho từng tham số tính trên sample thứ $w$:

$$\frac{\partial L^{(w)}}{\partial \Theta_{ij}} = \sum_{u=1}^{k} \frac{\partial L^{(w)}}{\partial \hat{y}_u^{(w)}} \times \sum_{v=1}^{k} \frac{\partial \hat{y}_v^{(w)}}{\partial z_i^{(w)}} \times \frac{\partial z_i^{(w)}}{d\Theta_{ij}} \tag{30}$$

3

Trong đó:

$$\frac{\partial z_i^{(w)}}{\partial \Theta_{ij}} = x_j^{(w)}, \tag{31}$$

$$\sum_{v=1}^{k} \frac{\partial \hat{y}_v^{(w)}}{\partial z_i^{(w)}} = \sum_{v=1}^{k} \frac{\partial}{\partial z_i^{(w)}} \left( \frac{e^{z_v^{(w)}}}{\sum_{j=1}^{k} e^{z_j^{(w)}}} \right) \tag{32}$$

$$\Leftrightarrow \sum_{v=1}^{k} \frac{\partial \hat{y}_v^{(w)}}{\partial z_i^{(w)}} = \frac{e^{z_i^{(w)}}}{\sum_{j=1}^{k} e^{z_j^{(w)}}} - \sum_{v=1}^{k} \frac{e^{z_i^{(w)}} e^{z_v^{(w)}}}{\left( \sum_{j=1}^{k} e^{z_j^{(w)}} \right)^2} \tag{33}$$

$$\Leftrightarrow \sum_{v=1}^{k} \frac{\partial \hat{y}_u^{(w)}}{\partial z_i^{(w)}} = \sum_{v=1}^{k} \hat{y}_v^{(w)} (\delta_v^{(w)} - \hat{y}_i^{(w)}), \tag{34}$$

$$\sum_{u=1}^{k} \frac{\partial L^{(w)}}{\partial \hat{y}_u^{(w)}} = -\sum_{u=1}^{k} \frac{\delta_u^{(w)}}{\hat{y}_u^{(w)}} \tag{35}$$

Do đó:

$$\frac{\partial L^{(w)}}{\partial \Theta_{ij}} = -\sum_{u=1}^{k} \frac{\delta_u^{(w)}}{\hat{y}_u^{(w)}} \times \sum_{v=1}^{k} \hat{y}_v^{(w)} (\delta_v^{(w)} - \hat{y}_i^{(w)}) \times x_j^{(w)} \tag{36}$$

$$\Leftrightarrow \frac{\partial L^{(w)}}{\partial \Theta_{ij}} = x_j^{(w)} (\hat{y}_i^{(w)} - \delta_i^{(w)}), \tag{37}$$

$$\frac{\partial L}{\partial \Theta_{ij}} = \sum_{w=1}^{m} \frac{\partial L^{(w)}}{\partial \Theta_{ij}} \tag{38}$$

$$\Leftrightarrow \frac{\partial L}{\partial \Theta_{ij}} = \frac{1}{m} \sum_{w=1}^{m} x_j^{(w)} (\hat{y}_i^{(w)} - \delta_i^{(w)}) \tag{39}$$

Vectorization:

$$\frac{\partial L}{\partial \Theta} = \frac{1}{m} \mathbf{X} (\hat{\mathbf{Y}} - \Delta)^T \tag{40}$$

- Cập nhật tham số:

$$\mathbf{\Theta} := \mathbf{\Theta} - \eta \times \frac{1}{m} \mathbf{X} (\hat{\mathbf{Y}} - \Delta)^T \tag{41}$$

**2.** Xây dựng công thức tính forward và backward (tính đạo hàm cho từng tham số) cho bài toán softmax regression dùng phương pháp dựa trên one-hot encoding.

(a) Stochastic gradient descent

(b) Batch gradient descent

**Lời giải**

(a) Stochastic gradient descent

4

Cài đặt:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \\ 1 \end{bmatrix} \tag{42}$$

$$\mathbf{y} = \begin{cases} y^{(1)}, \\ y^{(2)}, \\ \cdots, \\ y^{(k)} \end{cases} \tag{43}$$

Với:

$$\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \cdots \\ y^{(k)} \end{bmatrix} = \mathbf{I}_k = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_k \end{bmatrix} \tag{44}$$

i. Tính toán forward

$$\mathbf{z} = \mathbf{\Theta}^T \mathbf{x} \tag{45}$$

$$\hat{\mathbf{y}} = \frac{e^{\mathbf{z}}}{\sum e^{\mathbf{z}}} \tag{46}$$

Loss function:

$$L = -\mathbf{y} \log \hat{\mathbf{y}} \tag{47}$$

ii. Tính toán backward
   - Gradient cho từng tham số:

$$\frac{\partial L}{\partial \Theta_{ij}} = \sum_{u=1}^{k} \frac{\partial L}{\partial \hat{y}_u} \times \sum_{v=1}^{k} \frac{\partial \hat{y}_v}{\partial z_i} \times \frac{\partial z_i}{\partial \Theta_{ij}} \tag{48}$$

Trong đó:

$$\frac{\partial z_i}{\partial \Theta_{ij}} = x_j, \tag{49}$$

$$\sum_{v=1}^{k} \frac{\partial \hat{y}_v}{\partial z_i} = \sum_{v=1}^{k} \hat{y}_v (y_v - \hat{y}_i), \tag{50}$$

$$\sum_{u=1}^{k} \frac{\partial L}{\partial \hat{y}_u} = -\sum (\mathbf{y}^T \oslash \hat{\mathbf{y}}), \quad \forall \hat{y} \in \hat{\mathbf{y}} \neq 0 \tag{51}$$

Từ đó, ta tìm được gradient:

$$\frac{\partial L}{\partial \Theta_{ij}} = -\sum (\mathbf{y}^T \oslash \hat{\mathbf{y}}) \times \sum_{v=1}^{k} \hat{y}_v(y_v - \hat{y}_i) \times x_j \tag{52}$$

$$\Leftrightarrow \frac{\partial L}{\partial \Theta_{ij}} = x_j(\hat{y}_i - y_i) \tag{53}$$

Vectorization:

$$\frac{\partial L}{\partial \Theta} = \mathbf{x}(\hat{\mathbf{y}} - \mathbf{y}^T) \tag{54}$$

- Cập nhật tham số:

$$\mathbf{\Theta} := \mathbf{\Theta} - \eta \times \mathbf{x}(\hat{\mathbf{y}} - \mathbf{y}^T) \tag{55}$$

(b) Batch gradient descent
Cài đặt:

$$\mathbf{X} = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \cdots & x_1^{(m)} \\ x_2^{(1)} & x_2^{(2)} & \cdots & x_2^{(m)} \\ \vdots & \vdots & \ddots & \vdots \\ x_n^{(1)} & x_n^{(2)} & \cdots & x_n^{(m)} \\ 1 & 1 & \cdots & 1 \end{bmatrix} \tag{56}$$

$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}^{(1)} \\ \mathbf{y}^{(2)} \\ \cdots \\ \mathbf{y}^{(m)} \end{bmatrix} \tag{57}$$

$$\tag{58}$$

Với $u \in [1, m]$:

$$\mathbf{y}^{(u)} = \begin{cases} \mathbf{y}^{[1]}, \\ \mathbf{y}^{[2]}, \\ \vdots \\ \mathbf{y}^{[k]} \end{cases}$$

$$\begin{bmatrix} \mathbf{y}^{[1]} \\ \mathbf{y}^{[2]} \\ \vdots \\ \mathbf{y}^{[k]} \end{bmatrix} = \mathbf{I}_k = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

$$\hat{\mathbf{Y}} = \begin{bmatrix} \hat{y}_1^{(1)} & \hat{y}_1^{(2)} & \cdots & \hat{y}_1^{(m)} \\ \hat{y}_2^{(1)} & \hat{y}_2^{(2)} & \cdots & \hat{y}_2^{(m)} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{y}_k^{(1)} & \hat{y}_k^{(2)} & \cdots & \hat{y}_k^{(m)} \end{bmatrix} \tag{59}$$

i. Tính toán forward

$$\mathbf{Z} = \mathbf{\Theta}^T \mathbf{x} \tag{60}$$

$$\hat{\mathbf{y}}^{(u)} = \frac{e^{\mathbf{z}^{(u)}}}{\sum e^{\mathbf{z}^{(u)}}} \tag{61}$$

$$\Rightarrow \hat{\mathbf{Y}} = \begin{bmatrix} \hat{\mathbf{y}}^{(1)} & \hat{\mathbf{y}}^{(2)} & \cdots & \hat{\mathbf{y}}^{(u)} \cdots \hat{\mathbf{y}}^{(m)} \end{bmatrix} \tag{62}$$

Loss function:

$$L \;=\; -\frac{1}{m}\sum_{u=1}^{m}\sum_{i=1}^{k}\delta_i^{(u)}\log \hat{y}_i^{(u)} \tag{63}$$

$$\Rightarrow L \;=\; -\frac{1}{m}\sum_{u=1}^{m}\sum \mathbf{y}^{[u]}\log \hat{\mathbf{y}}^{(u)} \tag{64}$$

$$\Rightarrow L \;=\; -\frac{1}{m}\sum \left(\mathbf{y}^T\odot\log\hat{\mathbf{y}}\right) \tag{65}$$

ii. Tính toán backward
   - Gradient cho từng tham số tính trên sample thứ $w$:

$$\frac{\partial L^{(w)}}{\partial \Theta_{ij}} \;=\; \sum_{u=1}^{k}\frac{\partial L^{(w)}}{\partial \hat{y}_u^{(w)}}\times\sum_{v=1}^{k}\frac{\partial \hat{y}_v^{(w)}}{\partial z_i^{(w)}}\times\frac{\partial z_i^{(w)}}{d\Theta_{ij}} \tag{66}$$

Trong đó:

$$\frac{\partial z_i^{(w)}}{\partial \Theta_{ij}} \;=\; x_j^{(w)}, \tag{67}$$

$$\sum_{v=1}^{k}\frac{\partial \hat{y}_v^{(w)}}{\partial z_i^{(w)}} \;=\; \sum_{v=1}^{k}\hat{y}_v^{(w)}(y_v^{(w)}-\hat{y}_i^{(w)}), \tag{68}$$

$$\sum_{u=1}^{k}\frac{\partial L^{(w)}}{\partial \hat{y}_u^{(w)}} \;=\; -\sum_{u=1}^{k}\frac{y_u^{(w)}}{\hat{y}_u^{(w)}} \tag{69}$$

$$\Rightarrow\sum_{u=1}^{k}\frac{\partial L^{(w)}}{\partial \hat{y}_u^{(w)}} \;=\; \sum \mathbf{y}^{(w)}\oslash(\hat{\mathbf{y}}^{(w)})^T \tag{70}$$

Do đó:

$$\frac{\partial L^{(w)}}{\partial \Theta_{ij}} \;=\; -\sum \mathbf{y}^{(w)}\oslash(\hat{\mathbf{y}}^{(w)})^T\times\sum_{u=1}^{k}\frac{y_u^{(w)}}{\hat{y}_u^{(w)}}\times x_j^{(w)} \tag{71}$$

$$\Leftrightarrow\frac{\partial L^{(w)}}{\partial \Theta_{ij}} \;=\; x_j^{(w)}(\hat{y}_i^{(w)}-y_i^{(w)}), \tag{72}$$

$$\frac{\partial L}{\partial \Theta_{ij}} \;=\; \sum_{w=1}^{m}\frac{\partial L^{(w)}}{\partial \Theta_{ij}} \tag{73}$$

$$\Leftrightarrow\frac{\partial L}{\partial \Theta_{ij}} \;=\; \frac{1}{m}\sum_{w=1}^{m}x_j^{(w)}(\hat{y}_i^{(w)}-y_i^{(w)}) \tag{74}$$

Vectorization:

$$\frac{\partial L}{\partial \boldsymbol{\Theta}} \;=\; \frac{1}{m}\mathbf{X}(\hat{\mathbf{Y}}^T-\mathbf{Y}) \tag{75}$$

   - Cập nhật tham số:

$$\boldsymbol{\Theta} := \boldsymbol{\Theta}-\eta\times\frac{1}{m}\mathbf{X}(\hat{\mathbf{Y}}^T-\mathbf{Y}) \tag{76}$$

**3.** Cài đặt bài toán softmax regression cho data `iris_1D_2c.csv` bằng phương pháp dựa vào one-hot encoding.

1. Stochastic gradient descent

2. Batch gradient descent

**4.** Cài đặt bài toán softmax regression cho data `iris_full.csv` bằng phương pháp dựa vào one-hot encoding.

1. Stochastic gradient descent

2. Batch gradient descent

**Lời giải cho bài 3 & 4**

LINK NOTEBOOK