

$$X \in \mathbb{R}^{2 \times 3}$$

$$\begin{pmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \end{pmatrix}$$

$$W \in \mathbb{R}^{3 \times 5}$$

$$\begin{pmatrix} W_{11} & W_{12} & W_{13} & W_{14} & W_{15} \\ W_{21} & W_{22} & W_{23} & W_{24} & W_{25} \\ W_{31} & W_{32} & W_{33} & W_{34} & W_{35} \end{pmatrix}$$

$$y \in \mathbb{R}^{2 \times 5}$$

$$= \begin{pmatrix} y_{11} & y_{12} & y_{13} & y_{14} & y_{15} \\ y_{21} & y_{22} & y_{23} & y_{24} & y_{25} \end{pmatrix}$$

$$\begin{cases} y_{12} = X_{11}W_{12} + X_{12}W_{22} + X_{13}W_{32} \\ y_{22} = X_{21}W_{12} + X_{22}W_{22} + X_{23}W_{32} \end{cases}$$

$$\begin{cases} \frac{\partial y_{12}}{\partial W_{12}} = X_{11} \\ \frac{\partial y_{22}}{\partial W_{22}} = X_{21} \end{cases}$$

J_w has the same shape as w ($\mathbb{R}^{3 \times 5}$)

Assuming $L = f(y) \in \mathbb{R}$

$$J_{w_{12}} = \frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial y_{12}} \cdot \frac{\partial y_{12}}{\partial w_{12}} + \frac{\partial L}{\partial y_{22}} \cdot \frac{\partial y_{22}}{\partial w_{12}}$$

$\underbrace{\qquad\qquad\qquad}_{x_{11}} \qquad\qquad\qquad \underbrace{\qquad\qquad\qquad}_{x_{21}}$

$$\begin{cases} \frac{\partial y_{12}}{\partial w_{12}} = x_{11} \\ \frac{\partial y_{22}}{\partial w_{22}} = x_{21} \end{cases}$$

Assuming $J_y = \frac{\partial L}{\partial y} \in \mathbb{R}^{2 \times 5}$, the same shape as y

then $J_{w_{12}} = J_{y_{12}} \cdot x_{11} + J_{y_{22}} \cdot x_{21}$

a dot product between
and 2nd column of J_y
and 1st column of x

$$y = X W$$

2×5 2×3 3×5

$$J_W = X^T \cdot J_y$$

3×5 3×2 2×5

Similarly:

$$J_X = J_y \cdot W^T$$

2×3 2×5 5×3

$$y = X + b$$

2×5 2×5 2×5

$$J_X = J_y$$

$$J_b = J_y$$

In code, we refer

J_y : y-grad

J_W : w-grad

$$y = a \cdot b$$

$$\frac{\partial y}{\partial a} = b$$

$$L = f(y)$$

$$y = a + b$$

$$\frac{\partial L}{\partial a}$$

if b is 1×5 .

then J_b is $\sum_{i=1}^2 J_{y_i}$ ← row of y

$$y = a + b$$

$$\frac{\partial y}{\partial a} = 1$$

$$\Rightarrow \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial a} = \frac{\partial L}{\partial y}$$

For element-wise computation like ReLU

$$\text{ReLU}(x) = \max(0, x)$$

the gradient is element-wise mapping

Check out matrix cookbook for
gradient of other linear algebra operation

Tips:

the gradient of mapping $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$
has the shape of \mathbb{R}^m