

Linear Layer

$$X \in \mathbb{R}^{m \times n}$$

m = batch size

n = input dim

$$Y = Wx + b$$

$$W \in \mathbb{R}^{k \times n}, \quad b \in \mathbb{R}^k$$

In some places

k = batch size

also used

and m = output dim

k = output dim

n = Input dim

Forward Propagation

$$Y = Wx^T + b$$

$$X \in \mathbb{R}^{m \times n}$$

$$W \in \mathbb{R}^{k \times n}$$

$$b \in \mathbb{R}^k$$

$$Y \in \mathbb{R}^{k \times m}$$



We will take
transpose of Y to
get $\mathbb{R}^{m \times k}$

Suppose $m=1$ on batch size is 1
then $x \in \mathbb{R}^{1 \times n}$

$$y = w x^T + b$$

$$y_i = \sum_{j=1}^n w_{ij} x_j + b_i \quad \text{--- } ①$$

Now suppose for $m > 1$ on batch size is more than 1

$$y_{ij} = \sum_{l=1}^n w_{il} x_{lj} + b_i \quad \text{--- } ②$$

\Rightarrow For more than one input, we can repeat equation ①, same thing we do to find derivative.

Derivative

$$y = w x^T + b \quad w \in \mathbb{R}^{m \times n}$$

$$x \in \mathbb{R}^n$$

$$\frac{\partial L}{\partial w}, \frac{\partial L}{\partial b}, \frac{\partial L}{\partial x} \quad b \in \mathbb{R}^m$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w}$$

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

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$$

let's consider for only one case
means $x \in R^{1 \times n}$

Previous grad on $\frac{\partial L}{\partial y} \in R^{m \times 1}$ (for
(case $\frac{\partial L}{\partial y} \in R^{m \times n}$)

then

$$y = w x^T + b$$

by ①

$$y_i = \sum_{j=1}^n w_{ij} x_j + b_i$$

$$\frac{\partial Y_i}{\partial w} = \frac{1}{\partial w} \left(\sum_{j=1}^n w_{ij} x_j + b_i \right)$$

$$\Rightarrow \frac{1}{\partial w} \sum_{j=1}^n w_{ij} x_j + \frac{1}{\partial w} b_i$$

$$\Rightarrow \frac{1}{\partial w} \sum_{j=1}^n w_{ij} x_j$$

$$\left(\frac{\partial Y_i}{\partial w} \right)_{kl} = \begin{cases} x_k & , \text{ if } k=i \\ 0 & , \text{ otherwise} \end{cases}$$

————— (3) —————

$1 \leq k \leq m$
 $1 \leq l \leq n$

$$\Rightarrow \begin{bmatrix} w_{11} & \dots & w_{1n} \\ \vdots & & \vdots \\ w_m & \dots & w_{mn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} w_{11} x_1 + \dots + w_{1n} x_n \\ \vdots \\ w_{m1} x_1 + \dots + w_{mn} x_n \end{bmatrix}$$

$$\left[\frac{\partial}{\partial w_{11}} (\omega_{11}x_1 + \dots + \omega_{1n}x_n) \quad \dots \quad \frac{\partial}{\partial w_{m1}} (\omega_{11}x_1 + \dots + \omega_{1n}x_n) \right. \\ \vdots \\ \left. \frac{\partial}{\partial w_{m1}} (\omega_{11}x_1 + \dots + \omega_{1n}x_n) \quad \dots \quad \frac{\partial}{\partial w_{mn}} (\omega_{11}x_1 + \dots + \omega_{1n}x_n) \right]$$

$$\Rightarrow \begin{bmatrix} x_1 & \dots & x_n \\ \textcircled{O} \end{bmatrix} \longrightarrow \textcircled{u}$$

Now by equation ③ and ⑥ we can show
full derivative

$$\frac{\partial y}{\partial w} = \begin{bmatrix} x_1 & \dots & x_n \\ \textcircled{O} \end{bmatrix} \quad \dots \quad \begin{bmatrix} \textcircled{O} & \dots & x_n \end{bmatrix}$$

$$\frac{\partial y}{\partial w} \in \mathbb{R}^{m \times m \times n} \longrightarrow \textcircled{S}$$

Prev grad, which is gradient of next layer and it will be used in chain rule of derivative to find the grad of loss w.r.t to w then, to equation (5) we can write

$$\Rightarrow \text{grad} @ \frac{\partial y}{\partial w} \quad \left\{ \begin{array}{l} \text{here grad is not directly } \frac{\partial L}{\partial y}, \text{ because it} \\ \text{will be scalar} \end{array} \right.$$

here,

(Q) \Rightarrow matrix multiplication

$$\text{grad} \in R^{1 \times m}$$

$$\text{grad} @ \frac{\partial y}{\partial w} = \text{grad}^T @ x$$

Suppose

— (6)

$$\text{grad} = [g_1, \dots, g_m]$$

$$x = [x_1, \dots, x_n]$$

then equation ⑥ will be

$$\text{grad } @ \frac{\partial y}{\partial w} = \begin{bmatrix} g_1 \\ \vdots \\ g_m \end{bmatrix} [x_1 \dots x_n]$$

⑦

Now for batch size K the input x dim will be $x \in \mathbb{R}^{K \times n}$, and $\text{grad} \in \mathbb{R}^{K \times m}$. We have to find derivative by taking mean of every

$$\text{hence, grad} = \frac{\partial L}{\partial y} \in \mathbb{R}^{K \times m}$$

$$\frac{\partial L}{\partial w} = \frac{1}{K} \sum_{i=1}^K \text{grad}_{(i)}^T @ x$$

$$\frac{\partial L}{\partial w} = \frac{1}{K} (\text{grad}^T @ x)$$

$$X \in R^{K \times n}$$

$$\text{grad} \in R^{K \times m}$$

$$\frac{\partial L}{\partial w} \in R^{m \times n}$$

Now let's find $\frac{\partial L}{\partial b}$, suppose

$$Y \in R^{1 \times m}, X \in R^{1 \times n}, W \in R^{m \times n}, b \in R^{1 \times m}$$

$$Y = WX + b$$

By equation ①

$$Y_i = \sum_{j=1}^n w_{ij}x_j + b_i$$

$$\frac{\partial Y_i}{\partial b} = \frac{\partial}{\partial b} \left(\sum_{j=1}^n w_{ij}x_j + b_i \right)$$

$$\frac{\partial Y_i}{\partial b} = \frac{\partial}{\partial b} \sum_{j=1}^n w_{ij}x_j + \frac{\partial}{\partial b} b_i$$

$$\frac{\partial Y_i}{\partial b} = \frac{\partial}{\partial b} b_i$$

$$\left(\frac{\partial Y_i}{\partial b} \right)_j = \begin{cases} 1, & \text{if } i=j \\ 0, & \text{else} \end{cases}$$

$$1 \leq i \leq m$$

$$1 \leq j \leq m$$

————— ⑧

To understand equation 8 Suppose

$$b = [b_1, \dots, b_m]$$

$$Y = [Y_1, \dots, Y_m]$$

$$\frac{\partial L}{\partial Y} = [g_1, \dots, g_m]$$

then,

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} = \begin{bmatrix} \frac{\partial y_1}{\partial b_1} & \cdots & \frac{\partial y_m}{\partial b_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial b_m} & \cdots & \frac{\partial y_m}{\partial b_m} \end{bmatrix}$$

By equation ⑧

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} = \mathbf{I}_{m \times m}$$

So,

————— ⑨

$$\frac{\partial L}{\partial \mathbf{b}} = \frac{\partial L}{\partial \mathbf{y}} \frac{\partial \mathbf{y}}{\partial \mathbf{b}}$$

$$\frac{\partial L}{\partial \mathbf{b}} = \frac{\partial L}{\partial \mathbf{y}} \mathbf{I}$$

$$\boxed{\frac{\partial L}{\partial \mathbf{b}} = \frac{\partial L}{\partial \mathbf{y}}}$$

————— ⑩

Now for case when batch size = K
then,

$$X \in R^{K \times n}$$

$$Y \in R^{K \times m}$$

$$\frac{\partial L}{\partial Y} \in R^{K \times m}$$

We have to find equation (g) for each b_i , $1 \leq i \leq K$,

Now $Y \in R^{K \times m}$ then,

$$\left(\frac{\partial Y_i}{\partial b} \right)_{a,b} = \begin{cases} 1, & \text{if } a=b \\ 0, & \text{otherwise} \end{cases}$$

$$1 \leq a \leq m$$

$$1 \leq b \leq m$$

$$1 \leq i \leq K$$

(ii)

By equation ⑪, we can see we will get $k \times m$ Identity matrix.

To find

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

where,

$$\frac{\partial L}{\partial y} \in R^{K \times m}$$

$$\frac{\partial L}{\partial b} \in R^{K \times m \times m} \\ (\text{By eq ⑪})$$

Now each K vector in $\frac{\partial L}{\partial y}$ will multiply with respective matrices in $\frac{\partial L}{\partial b}$ (eq ⑪), in return we will get $K \times m$ matrices, $\left(\frac{\partial y}{\partial b}\right)_i = I$,

$1 \leq i \leq K$ (By eq(11)) then

$$\frac{\partial L}{\partial b} = \frac{1}{K} \sum_{i=1}^K \left(\frac{\partial L}{\partial y} \right)_i$$

prev
grad

Note:

- => when we are finding gradient for batch, consider each input in batch independent with each other. So while we pass grad to each layer, we can find different grad for different input in batch, and combine once we find for each input in batch by taking their avg or sum.
- => But instead of finding different

gradient and then take their avg.
we generally combine this process
by taking advantage of algebra
and make single matrix operation
 \Rightarrow so just see different independent
(collection of rows) input in batch,
instead of matrix.

Now let's find $\frac{\partial L}{\partial x}$, Suppose

$$x \in \mathbb{R}^{1 \times n}, y \in \mathbb{R}^{1 \times m}, b \in \mathbb{R}^{1 \times m}, \frac{\partial L}{\partial x} \in \mathbb{R}$$

$$y = wx + b$$

By equation ①

$$Y_i = \sum_{j=1}^n w_{ij}x_j + b_i$$

then

$$\frac{\partial Y_i}{\partial x} = \frac{\partial}{\partial x} \left(\sum_{j=1}^n w_{ij}x_j + b_i \right)$$

$$\frac{\partial Y_i}{\partial x} = \frac{\partial}{\partial x} \sum_{j=1}^n w_{ij}x_j + \frac{\partial}{\partial x} b_i$$

$$\frac{\partial Y_i}{\partial x} = \frac{\partial}{\partial x} \sum_{j=1}^n w_{ij}x_j$$

$$\frac{\partial Y_i}{\partial x_j} = w_{ij} \Rightarrow \boxed{\frac{\partial Y}{\partial x} = W}$$

$$\frac{\partial Y}{\partial x} \in \mathbb{R}^{m \times n}$$

Now to find the $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} @ \frac{\partial y}{\partial x}$

Suppose

$$\frac{\partial L}{\partial y} = [g_1, \dots, g_m]$$

$$\frac{\partial y}{\partial x} = W = \begin{bmatrix} w_{11} & \dots & w_{1n} \\ \vdots & \ddots & \vdots \\ w_{m1} & \dots & w_{mn} \end{bmatrix}$$

$$\frac{\partial L}{\partial x} = [g_1, \dots, g_n] \begin{bmatrix} w_{11} & \dots & w_{1n} \\ \vdots & \ddots & \vdots \\ w_{m1} & \dots & w_{mn} \end{bmatrix}$$

$$\frac{\partial L}{\partial x} = [g_1 w_{11} + \dots + g_m w_{m1}, \dots, g_1 w_{1n} + \dots + g_n w_{mn}]$$

— (12)

Now Suppose $X \in R^{K \times n}$, $Y \in R^{K \times m}$,

$b \in R^{1 \times m}$, $W \in R^{m \times n}$, $\frac{\partial L}{\partial Y} \in R^{1 \times m}$

$$\frac{\partial Y}{\partial X} \in R^{K \times n}$$

Now we have to find derivative
for each K batch, this grad
we will provide previous layer.

By equation ⑫ we can generalize
for K batch

$$\boxed{\frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y} @ \frac{\partial Y}{\partial X}}$$

where, $\frac{\partial L}{\partial X} \in R^{K \times n}$, $\frac{\partial L}{\partial Y} \in R^{1 \times m}$

$$\frac{\partial Y}{\partial X} \in R^{m \times n}$$

