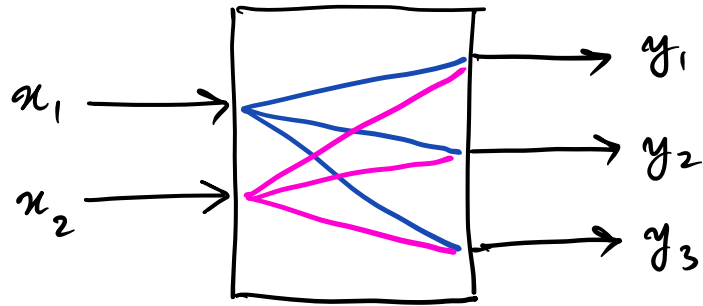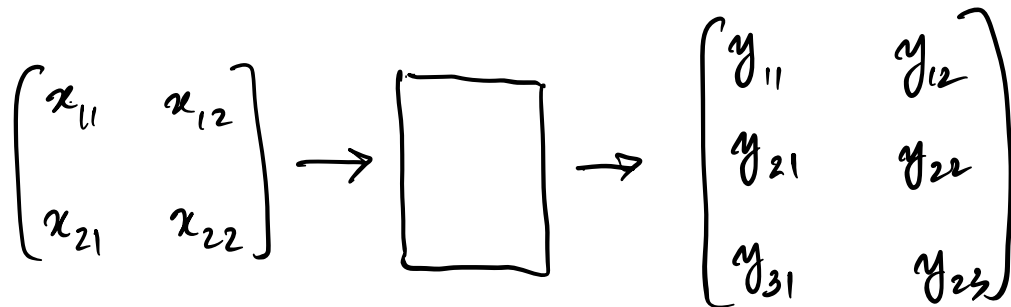# Linear layers



$$y_1 = W_{11} x_1 + W_{21} x_2$$

Re-write this as matrix-vector product

$$\begin{bmatrix} y_1 & y_2 & y_3 \end{bmatrix} = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \end{bmatrix}$$

Let's re-write this linear layer to handle mini-batches

$$\begin{pmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{pmatrix} \rightarrow \boxed{\phantom{XX}} \rightarrow \begin{pmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \\ y_{31} & y_{23} \end{pmatrix}$$

Batches

And in the matrix-vector form

$$\underset{\text{Batches}}{\downarrow}\begin{bmatrix} y_{11} & y_{21} & y_{31} \\ y_{12} & y_{22} & y_{32} \end{bmatrix} = \begin{bmatrix} x_{11} & x_{21} \\ x_{21} & x_{22} \end{bmatrix} \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{27} \end{bmatrix}$$
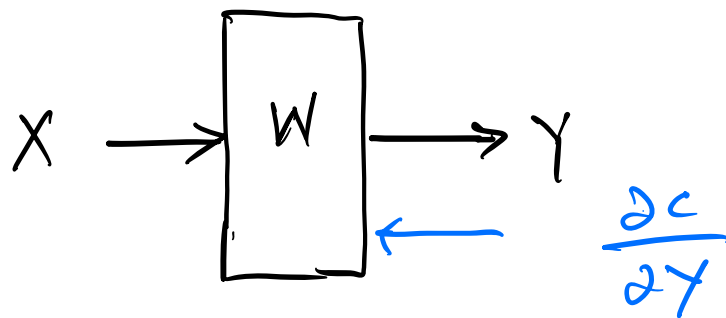
$$Y \qquad\qquad X \qquad\qquad W$$

For a linear layer, $\quad Y = XW \quad$ (Forward fn.)

In order to be able to backprop, we need to compute

$$\frac{\partial c}{\partial x} \quad \text{and} \quad \frac{\partial c}{\partial w}$$

Recall that we already have $\frac{\partial c}{\partial Y}$ (this was backpropagated)

In our example, $\frac{\partial c}{\partial Y}$ is the gradient whose size is the size of $Y$. $c$ is a scalar.

Therefore, $\frac{\partial c}{\partial Y} \in \mathbb{R}^{2 \times 3}$

$$\frac{\partial c}{\partial Y} = \begin{pmatrix} \frac{\partial c}{\partial y_{11}} & \frac{\partial c}{\partial y_{21}} & \frac{\partial c}{\partial y_{31}} \\ \\ \frac{\partial c}{\partial y_{12}} & \frac{\partial c}{\partial y_{22}} & \frac{\partial c}{\partial y_{32}} \end{pmatrix}$$

By applying chain-rule, we get

$$\frac{\partial c}{\partial x} = \frac{\partial c}{\partial Y} \cdot \frac{\partial Y}{\partial x}$$

Following from above

$$\frac{\partial c}{\partial x} = \begin{pmatrix} \frac{\partial c}{\partial x_{11}} & \frac{\partial c}{\partial x_{2,}} \\ \\ \frac{\partial c}{\partial x_{12}} & \frac{\partial c}{\partial x_{22}} \end{pmatrix}$$

Let's try to unpack $\dfrac{\partial Y}{\partial X}$. We begin by considering the one batch first.

$$\begin{bmatrix} y_{11} & y_{21} & y_{31} \end{bmatrix} = \begin{bmatrix} x_{11} & x_{21} \end{bmatrix} \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}$$

This suggests that $\dfrac{\partial Y}{\partial X}$ is actually a Jacobian.

$$\frac{\partial Y}{\partial X} = \begin{bmatrix} \dfrac{\partial y_{11}}{\partial x_{11}} & \dfrac{\partial y_{11}}{\partial x_{21}} \\[4mm] \dfrac{\partial y_{21}}{\partial x_{11}} & \dfrac{\partial y_{21}}{\partial x_{21}} \\[4mm] \dfrac{\partial y_{31}}{\partial x_{11}} & \dfrac{\partial y_{31}}{\partial x_{21}} \end{bmatrix} = \begin{bmatrix} w_{11} & w_{21} \\[2mm] w_{12} & w_{22} \\[2mm] w_{13} & w_{23} \end{bmatrix} = W^{T}$$

Now let's try to dig a little deeper. to see what is
$\dfrac{\partial C}{\partial x_{11}}$ ?

$$\frac{\partial c}{\partial x_{11}} = \frac{\partial c}{\partial y_{11}} \frac{\partial y_{11}}{\partial x_{11}} + \frac{\partial c}{\partial y_{21}} \frac{\partial y_{21}}{\partial x_{11}} + \frac{\partial c}{\partial y_{31}} \frac{\partial y_{31}}{\partial x_{11}}$$

$$= \begin{bmatrix} \dfrac{\partial c}{\partial y_{11}} & \dfrac{\partial c}{\partial y_{21}} & \dfrac{\partial c}{\partial y_{31}} \end{bmatrix} \begin{bmatrix} \dfrac{\partial y_{11}}{\partial x_{11}} \\[2ex] \dfrac{\partial y_{21}}{\partial x_{11}} \\[2ex] \dfrac{\partial y_{31}}{\partial x_{11}} \end{bmatrix}$$

We can re-write it for every dimension of a single sample.

$$\begin{bmatrix} \dfrac{\partial c}{\partial x_{11}} & \dfrac{\partial c}{\partial x_{21}} \end{bmatrix} = \begin{bmatrix} \dfrac{\partial c}{\partial y_{11}} & \dfrac{\partial c}{\partial y_{21}} & \dfrac{\partial c}{\partial y_{31}} \end{bmatrix} \begin{bmatrix} \dfrac{\partial y_{11}}{\partial x_{11}} & \dfrac{\partial y_{11}}{\partial x_{21}} \\[2ex] \dfrac{\partial y_{21}}{\partial x_{11}} & \dfrac{\partial y_{21}}{\partial x_{21}} \\[2ex] \dfrac{\partial y_{31}}{\partial x_{11}} & \dfrac{\partial y_{31}}{\partial x_{21}} \end{bmatrix}$$

$$\frac{\partial c}{\partial X} = \frac{\partial c}{\partial Y} W^T$$

This can be easily extended to the minibatch.

$$\begin{pmatrix} \frac{\partial c}{\partial x_{11}} & \frac{\partial c}{\partial x_{21}} \\ \\ \frac{\partial c}{\partial x_{12}} & \frac{\partial c}{\partial x_{22}} \end{pmatrix} = \begin{pmatrix} \frac{\partial c}{\partial y_{11}} & \frac{\partial c}{\partial y_{21}} & \frac{\partial c}{\partial y_{31}} \\ \\ \frac{\partial c}{\partial y_{12}} & \frac{\partial c}{\partial y_{22}} & \frac{\partial c}{\partial y_{32}} \end{pmatrix} W^T$$

We can similarly find $\frac{\partial c}{\partial W}$.

Useful properties of linear layers.

- Global context. Every output depends upon every input.
- Information mixin
- Frequently used as the last layer for many commonly used networks.

- Note that the above description ignores the activation functions.