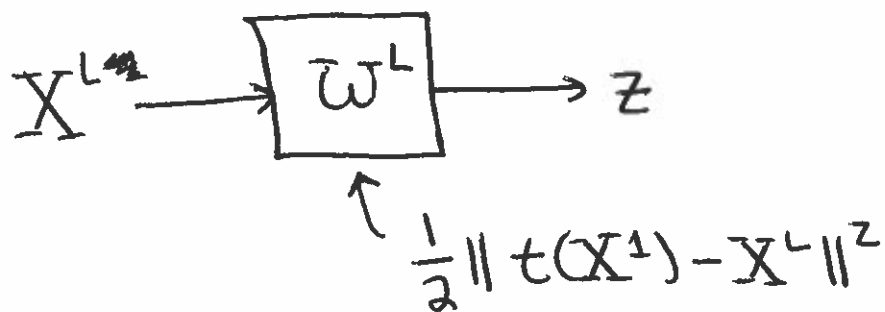Computing derivatives.

We need to compute $\frac{\partial z}{\partial w_i}$ for each $i$.

At the end of the network
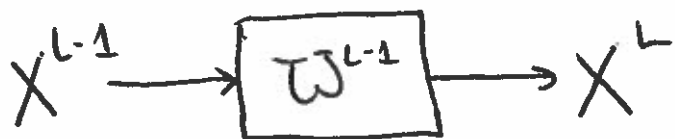
$$X^{L-1} \longrightarrow \boxed{w^L} \longrightarrow z$$

$$\uparrow \frac{1}{2}\| t(X^1) - X^L \|^2$$

it's easy to see

$$\frac{\partial z}{\partial w^L} = 0 \qquad \frac{\partial z}{\partial X^L} = \frac{\partial}{\partial X^L} \frac{1}{2} \left( t(X^1) - X^L \right).$$

$$\left( t(X^1) - X^L \right)$$

$$= X^L - t.$$

Now consider

$$X^{L-1} \longrightarrow \boxed{w^{L-1}} \longrightarrow X^L$$

We have (in principle)

$$\frac{\partial z}{\partial w^{L-1}} = \frac{\partial z}{\partial X^L} \cdot \frac{\partial X^L}{\partial w^{L-1}} \qquad \text{(chain rule)}$$

and

$$\frac{\partial z}{\partial X^{L-1}} = \frac{\partial z}{\partial X^L} \cdot \frac{\partial X^L}{\partial X^{L-1}} \quad \text{(chain rule)}$$

However, we need to keep track of the order (# of indices) and range $\left(\begin{smallmatrix} \text{values} \\ \text{for index} \end{smallmatrix}\right)$ of ~~each~~ the tensors $X^L, X^{L-1}$ to make the multiplications above make sense.

So we write

$$\text{vec}(X^L) \text{ as a function of } \text{vec}(X^{L-1}),$$
$$\text{vec}(W^{L-1})$$

and then

$$\frac{\partial z}{\partial \text{vec}(W^{L-1})^T} = \underset{\displaystyle\downarrow}{\overset{\text{vectors}}{\frac{\partial z}{\partial \text{vec}(X^L)^T}}} \cdot \underset{\displaystyle\downarrow}{\overset{\text{matrices}}{\frac{\partial \text{vec}(X^L)}{\partial \text{vec}(W^{L-1})^T}}}$$

$$\frac{\partial z}{\partial \text{vec}(X^{L-1})^T} = \frac{\partial z}{\partial \text{vec}(X^L)^T} \cdot \frac{\partial \text{vec}(X^L)}{\partial \text{vec}(X^{L-1})^T}$$

makes sense.

Thus is if we can compute $\xi$ for each

layer $X^i \longrightarrow \boxed{W^i} \longrightarrow X^{i+1}$

$$\frac{\partial vec(X^{i+1})}{\partial vec(W^i)^T}$$  (derivatives of output tensor w.r.t. parameters)

and

$$\frac{\partial vec(X^{i+1})}{\partial vec(X^i)^T}$$  (derivatives of output tensor w.r.t. input)

we can compute

$$\frac{\partial z}{\partial vec(W^i)^T}$$

by working our way <u>backwards</u> from $W^L$ to $W^1$. This is called back propagation.
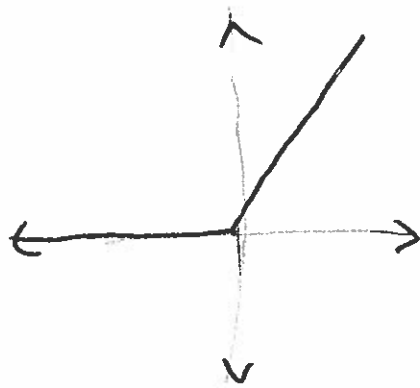
Now we need to know:

What kind of layers do we want?

What are their derivatives?

~~ReLu~~

ReLU or Ramp layers.

Consider the function $r(x) = \max(0, x)$.



$$r'(x) = \begin{cases} 1, & x > 0 \\ \text{undefined,} & x = 0 \\ 0, & x < 0 \end{cases}$$

a ramp layer applies $r(x)$ to every entry in $\underline{X}^i$, returning an output $\underline{X}^{i+1}$ of the same shape.
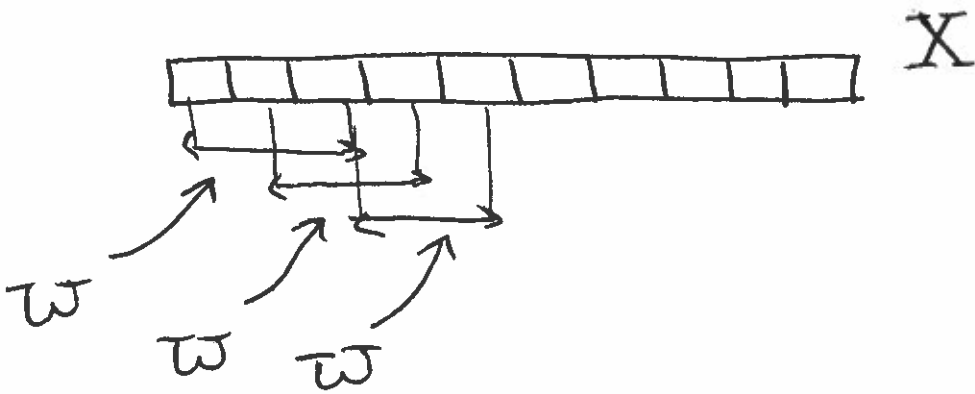
Convolution Layers.

Convolution is an operation ~~with~~ which combines two tensors in a ~~&~~ quadratic way.

Example. $X$ is a vector with 10 entries. $W$ is a vector with 3 entries. $X * W$ is a vector with 8 entries defined by

$$(X * W)_i = \underbrace{\sum_{j=1}^{3} W_j X_{i+j-1}}_{\substack{\text{dot product} \\ \text{of } W_{\cancel{3}} \text{ with} \\ \underline{\text{part}} \text{ of } \overline{X}}}$$

The $W$ tensor is called a <u>convolution</u>
<u>kernel</u>. The ~~convolved~~ <u>convolution</u>
<u>product</u> $X * W$ is large when $\underline{X}$
contains <u>features</u> detected by $W$.

Example. Suppose $X_i = f(ih)$ for
some differentiable function $f$ and
small $h$, and $W = (1, -2, 1)$.
Then

$$(W * \underline{X})_i = \underline{X}_i - 2\underline{X}_{i+1} + \underline{X}_{i+2}$$

$$= f(ih) - 2f((i+1)h) + f((i+2)h).$$

or if $(i+1)h = x_0$,

$$= f(x_0 - h) - 2f(x_0) + f(x_0 + h).$$

To under this term, we expand $f(x_0 \pm h)$ using Taylor's theorem

$$= f(x_0) - hf'(x_0) + \frac{h^2}{2}f''(x_0) - \frac{h^3}{6}f'''(x_0) + \ldots$$
$$- 2f(x_0)$$
$$f(x_0) + hf'(x_0) + \frac{h^2}{2}f''(x_0) + \frac{h^3}{6}f'''(x_0) + \ldots$$

$$= h^2 f''(x_0) + O(h^4)$$

This is called the Sobel Kernel and it (roughly) detects edges in an image.

Convolution is <u>linear</u> in $\underline{X}$ and <u>linear</u> in $W$, so the derivatives

$$\frac{\partial(\underline{X} * W)}{\partial \underline{X}}, \quad \frac{\partial(\underline{X} * W)}{\partial W}$$

are easy to compute.

Pooling layer.

The pooling layer replaces $\overset{\text{disjoint}}{\wedge}$ regions of the input tensor with their max ~~of~~ or average. There are no parameters.