

In the case of an image, we have a 3D array, but instead of considering this finite scenario, we will look at \mathbb{R}^3 or in other words "an infinite image". Almost the same calculations will be involved, but with more labels in the finite case.

Assume we have a function on \mathbb{R}^3 , $f: \mathbb{R}^3 \rightarrow \mathbb{R}^C$. It could be the output of a ^{neural network} that assigns to the initial image, its feature map. So, C denotes the channel feature map.

Consider the space

$$L^2(\mathbb{R}^3, \mathbb{R}^C) := \left\{ F: \mathbb{R}^3 \rightarrow \mathbb{R}^C \mid \int_{\mathbb{R}^3} dx \|F(x)\|^2 < \infty \right\}$$

of square integrable functions. On it, we can define a translation operator $T_t: L^2 \rightarrow L^2$,

$$(T_t f)(x) = f(x - t)$$

This is in other words, a representation of

$$(\mathbb{R}^3, +) \longrightarrow GL(L^2(\mathbb{R}^3, \mathbb{R}^C)).$$

We will next look at an ansatz of how the operation of a neural network could look like:

$$(\mathcal{I}_K f)(x) := \int_{\mathbb{R}^3} dy K(x, y) f(y)$$

$$\mathcal{I}_K : L^2(\mathbb{R}^3, \mathbb{C}^{in}) \rightarrow L^2(\mathbb{R}^3, \mathbb{C}^{out}), \quad K: \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{R}^{C^{out} \times C^{in}}$$

Where \mathbb{C}^{in} , \mathbb{C}^{out} can represent the input and output feature maps of the neural network.

First symmetry (or equivariance) is translational. If our \mathcal{I}_K satisfies this, then

$$\mathcal{I}_K(T_t f) = T_t(\mathcal{I}_K f), \quad \forall f \in L^2(\mathbb{R}^3).$$

In other words, if the input is translated, then the feature maps translate by the same factor.

$$\begin{aligned}
\mathbb{I}_k(T_t t)^{(x)} &= \int_{\mathbb{R}^3} dy K(x, y) f(y-t) \\
&= \int_{\mathbb{R}^3} d\tilde{y} K(x, \tilde{y}+t) f(\tilde{y}) \\
&= (T_t (\mathbb{I}_k t))(x) = (\mathbb{I}_k t)(x-t) \\
&= \int_{\mathbb{R}^3} dy K(x-t, y) f(y)
\end{aligned}$$

The two are actually equal $\Leftrightarrow K(x-t, y) = K(x, y+t), (\forall) y \in \mathbb{R}^3$.

If we change y to $y-t \Rightarrow K(x, y) = K(x-t, y-t)$

So, the kernel has a "spatial weight sharing", the weights at (x, y) are the same weights at $(x-t, y-t)$.

If we choose $t = y \Rightarrow K(x, y) = K(x-y, 0) \Rightarrow$ the kernel depends only on $x-y$. So,

$$K: \mathbb{R}^3 \rightarrow \mathbb{R}^{c_{in} \times c_{out}}$$

If we have nonlinearities that suppose they are local, meaning $(\nabla f)(x) := \nabla_x (f(x))$ a different function at each $x \in \mathbb{R}^3$. If we require symmetry:

$$(\nabla (T_t f))(x) = (T_t (\nabla f))(x)$$

$$\text{Hence, } (\nabla (T_t f))(x) = \nabla_x (T_t f(x))$$

$$= \nabla_x (f(x-t)) = (T_t (\nabla f))(x) = (\nabla f)(x-t)$$

$$= \nabla_{x-t} (f(x-t))$$

$\Leftrightarrow \nabla_x = \nabla_{x-t} \Rightarrow \nabla_x = \nabla_0 := \nabla$. Therefore, the nonlinearity is a constant, independent of the position x .

Let's have a look at another crucial operation in convolutional networks i.e. local max pooling (LMP).

$$\text{LMP} : L^2(\mathbb{R}^3, \mathbb{R}^C) \rightarrow L^2(\mathbb{R}^3, \mathbb{R}^C),$$

$$f \mapsto \max_{y \in R_x} f(y), \text{ where } R_x \text{ is a region}$$

where the operation is performed.

Symmetry again, requires:

$$(\text{LMP}(T_t f))(x) = (T_t(\text{LMP} f))(x)$$

$$\Leftrightarrow \max_{y \in R_x} (T_t f)(y) = \max_{y \in R_x} f(y-t) = \max_{y \in R_{x-t}} f(y)$$

$$= (T_t(\text{LMP} f))(x) = (\text{LMP} f)(x-t) = \max_{y \in R_{x-t}} f(y)$$

$\Leftrightarrow R_x - t = R_{x-t} \Leftrightarrow$ the pooling windows are shared (have the same size in the context of CNN's.).