

why MLP failed?

- Locality destruction
- No parameter sharing
- Parameter Explosion

How we solve MLP's Problem?

→ we apply Mathematical Constraints

Mathematically

a) Locality destruction

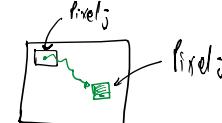
each MLP layer performs,

$$y = \omega x + b \quad \omega \in \mathbb{R}^{H \times W \times C}; \quad b \in \mathbb{R}^{d \times (H \times W \times C)}$$

$\Rightarrow y_i = \sum_{j=0}^N w_j x_j + b_j$ This is global linear operator over the image

\Rightarrow right now pixel; meaning

is influencing pixel;
and they are not even close. (neighbour).



Notes

▷ Right now each w_j is contributing for the prediction of y_i → images need to load pixels to only take part strongly

⇒ MLP is permutation-sensitive but Geometrically blind.

b) No Parameter sharing

$$y_i = \sum_{j=0}^N w_j x_j + b_j$$

Assume:- two pixel x_p and x_q represents same local pattern at two different location

Because right now there is no parameter sharing.

In General, $[w_{x_p} \neq w_{x_q}]$ ← model will have to re-learn this pattern at different location.

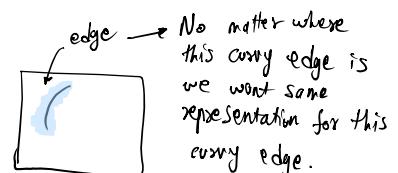
← Computation is wasted.

Notes

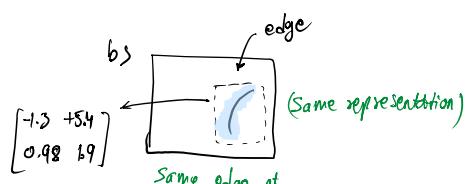
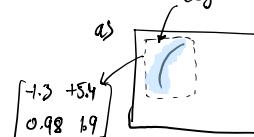
In current configuration,

▷ Model must relearn the same feature at every different position

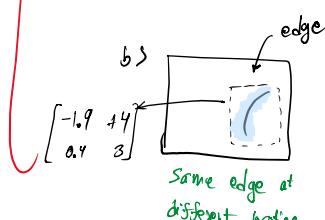
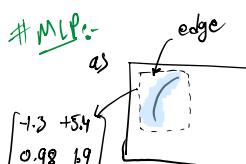
⇒ Translation of a pattern produce different representation.



what we want (Convolution) (weight sharing)



The same local pattern at two different location is processed by different linear maps.

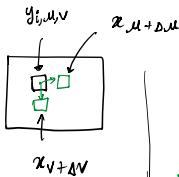


Parameter Explosion

MLP → Parameter Count for Single Layer in MLP,
 $O = (w \times h \times c \times d)$

Minimal Constraint on Linear Operator

↳ locality Dependent
 $y_{i,u,v}$ must be dependent only on neighbor pixels.
 $x_{u+\Delta u, v+\Delta v}$



Notation

$A_{i,u,v,c,u',v'} =$ weight connecting input signal (c,u',v') to output signal (i,u,v) .

⇒ $A_{i,u,v,c,u',v'} = 0$ unless $|u'-u| \leq K$ and $|v'-v| \leq K$.

⇒ Pixel that are not in a range of " $K \times K$ " will not contribute for predicting output signal.

⇒ Define offset, $\Delta u = u' - u$, $\Delta v = v' - v$

So operator becomes,

eqn ①

$$y_{i,u,v} = \sum_{c=1}^C \sum_{\Delta u, \Delta v} A_{i,u,v,c,u+\Delta u, v+\Delta v} \cdot x_{c,u+\Delta u, v+\Delta v}$$

locality Achieved
 But weights are still dependent upon absolute position (u,v) .

2) Weight sharing

↳ we want to store the same weight across a range of " $K \times K$ ".

→ offset.

4.

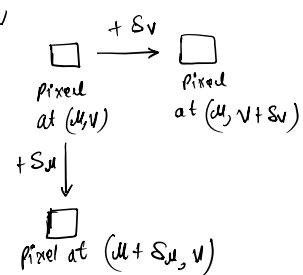
How do we shift the image??

→ we define " T_δ " → shift the image.

$$\Rightarrow S = \{S_u, S_v\}$$

S_u ↑
 up/down shift
 S_v ↑
 Right-left shift

We want,



equivalence

$x \rightarrow$ input space, $y \rightarrow$ output space.

g : group acting on both input and output spaces

$$T: X \rightarrow Y$$

$$T(gx) = g T(x); \quad \text{for } g$$

Then, T is equivariant to g

for image. $G \rightarrow g \rightarrow T_\delta \rightarrow$ Spatial shift.

$$\Rightarrow T(T_\delta x) = T_\delta T(x)$$

This will force weight sharing

$$(T_\delta x)_{c,u,v} = x_{c,u-S_u, v-S_v}$$

pixel at (u,v) is just a shifted version of pixel at $(u-S_u, v-S_v)$

what we have now,

$$y_{i,u,v} = \sum_{c=1}^C \sum_{\Delta u, \Delta v} A_{i,u,v,c,\Delta u, \Delta v} \cdot x_{c,u+\Delta u, v+\Delta v}$$

pixel at (u, v) is just
shifted version of pixel
at $((u-s_u), (v-s_v))$

shift input by s

$$\Rightarrow T(\mathcal{C}_s x)_{i,u,v} = \sum_{c=1}^C \sum_{\Delta u, \Delta v} A_{i,u,v,c,\Delta u, \Delta v} \cdot x_{c,u+s_u, v+s_v} \rightarrow \text{LHS.}$$

$$\Rightarrow T_s(T(x))_{i,u,v} = \sum_{c=1}^C \sum_{\Delta u, \Delta v} A_{i,u-s_u, v-s_v, c, \Delta u, \Delta v} \cdot x_{c,u+s_u, v+s_v} \rightarrow \text{RHS}$$

But, $T(\mathcal{C}_s x) = T_s(T(x))$

$$\sum_{c=1}^C \sum_{\Delta u, \Delta v} A_{i,u,v,c,\Delta u, \Delta v} \cdot x_{c,u+s_u, v+s_v} = \sum_{c=1}^C \sum_{\Delta u, \Delta v} A_{i,u-s_u, v-s_v, c, \Delta u, \Delta v} \cdot x_{c,u+s_u, v+s_v}$$

original position \downarrow This solution enables weight sharing.

$$\Rightarrow A_{i,u,v,c,\Delta u, \Delta v} = A_{i,u-s_u, v-s_v, c, \Delta u, \Delta v}$$

\downarrow shifted position. Because weight at (u, v) is used at $(u-s_u, v-s_v)$.

\Rightarrow Independent of u, v .
 \Rightarrow weight don't depends on absolute position

Meaning

$$A_{i,u,v,c,\Delta u, \Delta v} = A_{i,u-s_u, v-s_v, c, \Delta u, \Delta v}$$

\Rightarrow The weight used at output position (u, v) is exactly the same as the weight used at output position $(u-s_u, v-s_v)$ as long as the relative offset $(\Delta u, \Delta v)$ is the same.

so if,

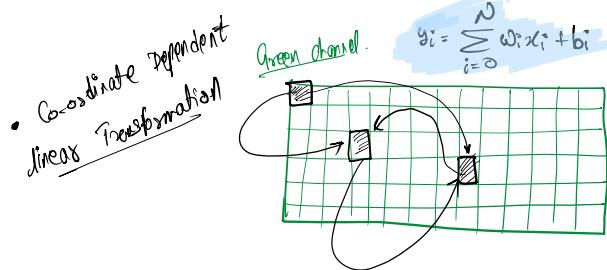
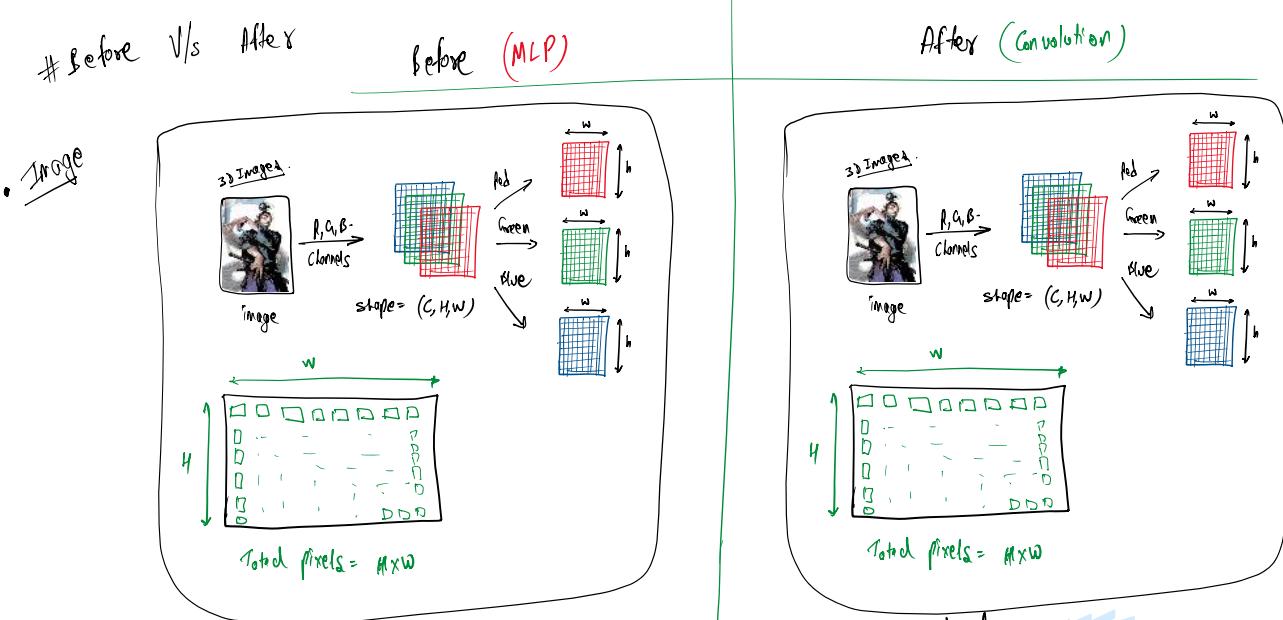
$$A_{i,u,v,c,\Delta u, \Delta v} = A_{i,u-s_u, v-s_v, c, \Delta u, \Delta v} - \text{eq } 10 \quad (2)$$

\Rightarrow it mustn't depends on (u, v) ; Because it satisfies eq 10 (2)

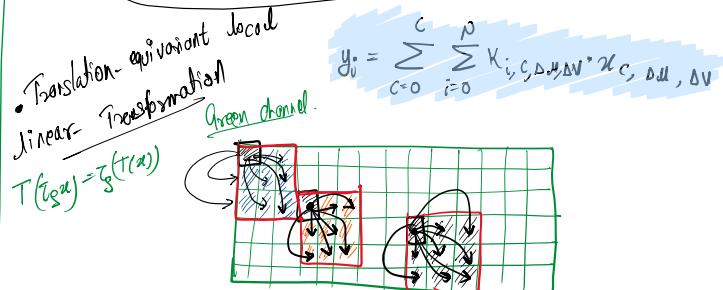
\Rightarrow rename it as,

Convolution Kernel,

$$A_{i,u,v,c,\Delta u, \Delta v} = K_{i,c,\Delta u, \Delta v}$$



→ All pixel will effect each
→ each one have different weight.



→ Pixel only effects the pixels within the red '3x3' window size (Kernel size)
→ only local pixels influence the output
→ distant pixel have zero contribution

$$A_{i,\mu,\nu, c, \mu', \nu'} = 0; \text{ unless } |\mu' - \mu| \leq K \text{ and } |\nu' - \nu| \leq K$$

→ weights are shared across spatial location
not across offsets.

Conclusion.

Convolution.

$$y_{i,\mu,\nu} = \sum_{c=1}^C \sum_{\Delta\mu=-K'}^{K'} \sum_{\Delta\nu=-K'}^{K'} k_{i,c,\Delta\mu,\Delta\nu} \cdot x_{c,\mu+\Delta\mu, \nu+\Delta\nu} + b_i$$

or,

$$y_{i,\mu,\nu} = \sum_{c,\Delta\mu,\Delta\nu} k_{i,c,\Delta\mu,\Delta\nu} \cdot x_{c,\mu+\Delta\mu, \nu+\Delta\nu} + b_i$$

MLP.

$$y_i = \sum_{j=1}^{CHW} w_{i,j} \cdot x_j + b_i$$

Constrained. → locality
→ weight sharing

⇒ A Convolution is just a linear layer

\Rightarrow A Convolution is just a linear layer
whose weight matrix has been constrained
to be sparse (local) and structured (shared across spatial shift)