

Geometric Learning

Question: Can we reduce the complexity of our hypothesis class without increasing the approximation error?

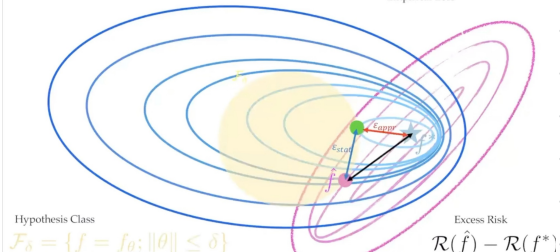
Answer: Using geometric priors (e.g. symmetry and scale separations), we can decrease the complexity of the hypothesis class while keeping the approximation error constant.

Recap: Error Analysis

$$\leq \epsilon_{\text{opt}} + 2 \sup_{f \in \mathcal{F}_h} |R(f) - \hat{R}(f)| + \epsilon_{\text{appr}} = \epsilon_{\text{opt}} + \epsilon_{\text{stat}} + \epsilon_{\text{appr}}$$

$$\hat{f} = \arg \min_{f \in \mathcal{F}_h} \hat{R}(f) \quad \mathcal{R}(f) = \mathbb{E}[\ell(f(x), y)] \quad \hat{\mathcal{R}}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

ERM Population Loss Empirical Loss



How to reduce ϵ_{opt} ? : Use a better optimization

ϵ_{appr} ? : Use a larger hypothesis class \mathcal{F}_h

ϵ_{stat} ? : # of data samples,
Complexity of the hypothesis class, \mathcal{F}_h

$$\epsilon_{\text{stat}} = \mathcal{R}(\hat{f}) - \hat{\mathcal{R}}(\hat{f}) \leq c \cdot \frac{\text{Complexity}}{n^{\frac{1}{d}}}$$

Error Analysis

: Let's say we found $\hat{f} \in \mathcal{F}_h$ to solve the problem. How good is it?

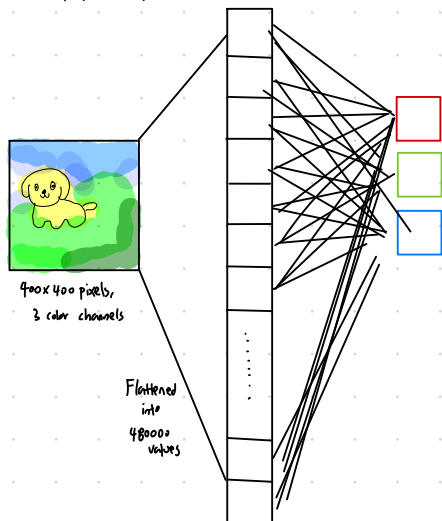
Mathematically, how small is the population error? : $\mathcal{R}(\hat{f}) - \inf_{f \in \mathcal{F}_h} \mathcal{R}(f)$

Statistical (ϵ_{stat}): Gap between $\mathcal{R}(\hat{f})$ and $\hat{\mathcal{R}}(\hat{f})$

Optimization (ϵ_{opt}): How far is \hat{f} from the best $f \in \mathcal{F}_h$

Approximation (ϵ_{appr}): How good is \mathcal{F}_h ? How far is the best $f \in \mathcal{F}_h$ from f^* ?

Where MLP (Multi-Layer Perceptrons) fail about?



However,



\neq



MLP model cannot distinguish transformed (shifted, rotated, reflected) data, even data is symmetry.

Similarly, changing the indices of the graph won't change the underlying graph.

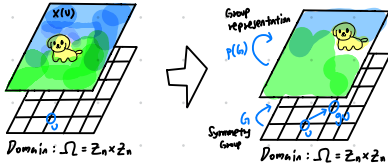
Geometric Priors

Geometric prior: The input is a signal defined on a domain with a geometric structure.

Signals $X(\Omega) := \{x | x: \Omega \rightarrow \mathbb{C}\}$

$$C = \mathbb{R}^3$$

$$P(g)x(u) = x(g^{-1}u)$$



Definition of Group

A group is a set G along with a binary operation $\circ: G \times G \rightarrow G$

① Associativity: $(g \circ h) \circ f = g \circ (h \circ f) \quad \forall g, h, f \in G$

② Identity: $I \circ g = g \circ I = g \quad \forall g \in G$

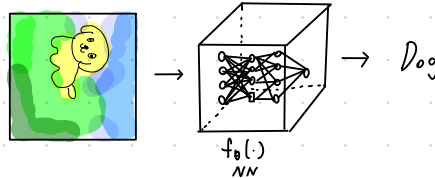
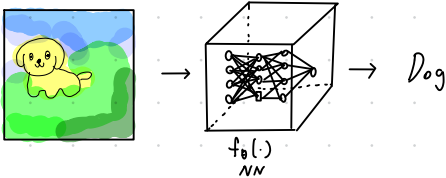
③ Inverse: $g \circ g^{-1} = I = g^{-1} \circ g$

④ Closure: $g \circ h \in G, \quad \forall g, h \in G$

For brevity we often write gh instead of $g \circ h$

Geometric Priors: G -invariance

$f_\theta(P(g)x) = f_\theta(x)$, where $P(g)x$: transformed signal



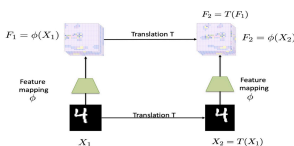
Geometric Priors: G -equivariance

$$f_\theta(P(g)x) = P(g)f_\theta(x)$$

How can we build equivariant functions?

Meaning of equivariant function: $f_\theta(P(g)x) = P(g)f_\theta(x)$

Pictorial depiction



$$F_1 = \phi(X_1) = f_\theta(X_1)$$

$$T = P(g)$$

$$X_2 = P(g)X_1 = T(X_1)$$

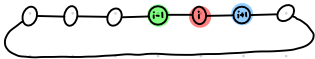
Then, as you can see,

$$T(F_1) = \phi(X_2)$$

$\Rightarrow P(g)f_\theta(x) = f_\theta(P(g)x)$, hence it's equivariant.

From Christian Wolf

Grid: ring graph



local aggregation function

$$f(x_i) = \phi(x_i, \{x_{i-1}, x_{i+1}\}) = \phi(x_{i-1}, x_i, x_{i+1})$$

Linear local aggregation function

$$f(x_i) = ax_{i-1} + bx_i + cx_{i+1}$$

Convolution \rightarrow vector of parameters θ

$$f(X) = \begin{bmatrix} b & c & 0 & 0 & a \\ a & b & c & 0 & 0 \\ & & \dots & & \\ 0 & 0 & a & b & c \\ c & 0 & 0 & a & b \end{bmatrix} \begin{bmatrix} X \\ X \\ X \\ X \\ X \end{bmatrix}$$

Circulant matrix
= convolution

ex) image:

x_1	x_2	x_3	x_4
x_5	x_6	x_7	x_8
x_9	x_{10}	x_{11}	x_{12}
x_{13}	x_{14}	x_{15}	x_{16}

Convolution =

k_1	k_2	k_3	0	k_4	k_5	k_6	0	k_7	k_8	k_9	0	0	0	0	0
0	k_1	k_2	k_3	0	k_4	k_5	k_6	0	k_7	k_8	k_9	0	0	0	0
0	0	0	0	k_1	k_2	k_3	0	k_4	k_5	k_6	0	k_7	k_8	k_9	0
0	0	0	0	0	k_1	k_2	k_3	0	k_4	k_5	k_6	0	k_7	k_8	k_9

Deriving Convolution from Symmetry

$$\begin{bmatrix} b & c & 0 & 0 & a \\ a & b & c & 0 & 0 \\ 0 & a & b & c & 0 \\ 0 & 0 & a & b & c \\ c & 0 & 0 & a & b \end{bmatrix} \begin{bmatrix} y & z & 0 & 0 & x \\ x & y & z & 0 & 0 \\ 0 & x & y & z & 0 \\ 0 & 0 & x & y & z \\ z & 0 & 0 & x & y \end{bmatrix} = \begin{bmatrix} y & z & 0 & 0 & x \\ x & y & z & 0 & 0 \\ 0 & x & y & z & 0 \\ 0 & 0 & x & y & z \\ z & 0 & 0 & x & y \end{bmatrix} \begin{bmatrix} b & c & 0 & 0 & a \\ a & b & c & 0 & 0 \\ 0 & a & b & c & 0 \\ 0 & 0 & a & b & c \\ c & 0 & 0 & a & b \end{bmatrix}$$

Circulant matrices commute

Ex)

Convolution \rightarrow shifting

$$\begin{bmatrix} b & c & 0 & 0 & a \\ a & b & c & 0 & 0 \\ 0 & a & b & c & 0 \\ 0 & 0 & a & b & c \\ c & 0 & 0 & a & b \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} b & c & 0 & 0 & a \\ a & b & c & 0 & 0 \\ 0 & a & b & c & 0 \\ 0 & 0 & a & b & c \\ c & 0 & 0 & a & b \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} a & b & c & 0 & 0 \\ 0 & a & b & c & 0 \\ 0 & 0 & a & b & c \\ c & 0 & 0 & a & b \\ b & c & 0 & 0 & a \end{bmatrix} = \begin{bmatrix} a & b & c & 0 & 0 \\ 0 & a & b & c & 0 \\ 0 & 0 & a & b & c \\ c & 0 & 0 & a & b \\ b & c & 0 & 0 & a \end{bmatrix}$$

Hence, $CS = SC$ iff Convolution is circulant matrix.

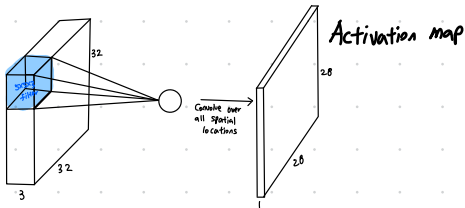
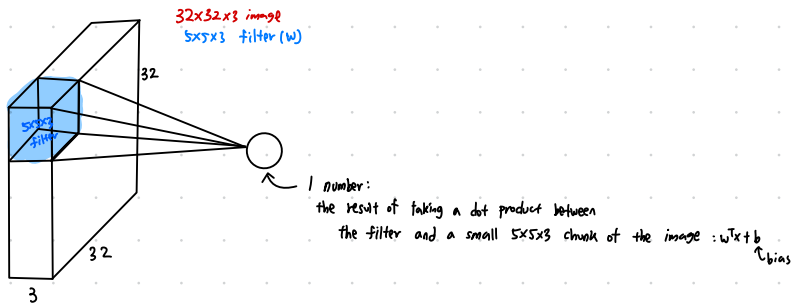
Recap: Geometric Priors: G -equivariance

$$f_{\theta}(P(g)x) = P(g)f_{\theta}(x)$$

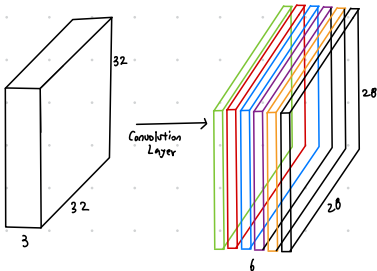
Therefore, Convolution matrix (circulant matrix) is shift-equivariant.

Convolution \Leftrightarrow shift-equivariant.

Convolution Operator

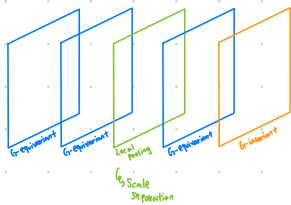


For example, if we had 6 5×5 filters, we will get 6 separate activation maps:

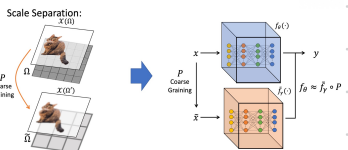


We stack these up to get a new image of size $28 \times 28 \times 6$

Geometric Deep Learning - A Blueprint

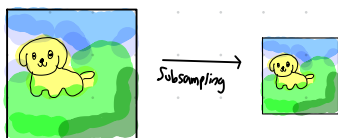


Scale Separation:



Why pooling?

: Subsampling pixels will not change the object



We can subsample the pixels to make image smaller, fewer parameters to characterize the image

Max Pooling

