

Lecture 11 : Learning Under Invariances and Geometric Priors

Recap last week: Benefits of depth in NNs in terms of approximation.

Today: . Look for other sources of regularity ; how to incorporate them into learning model.

- Two key principles: (i) symmetries / invariances
- (ii) multiscale separation.

Symmetries and Geometric Stability

→ So far, we have measured the regularity / complexity of a target function $f: X \rightarrow \mathbb{R}$ with norms $\gamma(f)$ of the form

$$\gamma(f) = \int |\hat{f}(\vec{z})|^2 (1 + \|\vec{z}\|^2)^s d\vec{z} \quad [\text{Sobolev norm}]$$

$$\gamma(f) = \int |\hat{f}(\vec{z})| \|\vec{z}\|^k d\vec{z} \quad [\text{Barron Norm}]$$

$$\gamma(f) = \inf \left\{ \|g\|_{L^2(\Theta, \tau)}^2 ; f(x) = \int g(\theta) \pi(x, \theta) \tau(d\theta) \right\} \quad [\text{RKHS norm}]$$

→ So far, we have treated X as a "more" high-dim vector space — invariant / agnostic to how we parametrize the input.!

→ In most situations of interest, the input space X has much more structure:

$$X = \{x : \Omega \rightarrow \mathcal{C}\}$$

↓
physical domain

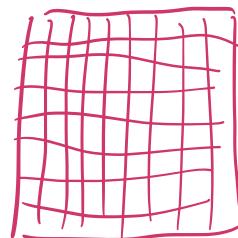
↓
feature channels

ex:

images

$$\Omega = [0, 1]^2$$

discretised version



$$\Omega = [1, n] \times [1, n]$$

\mathcal{C} : RGB space of colors.

speech/music

$$\Omega = \mathbb{R}$$

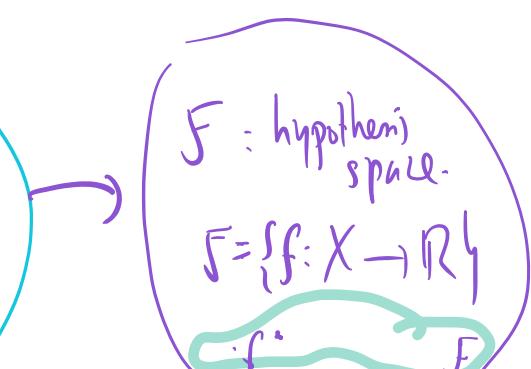
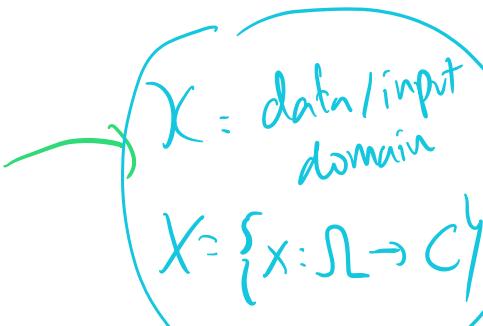
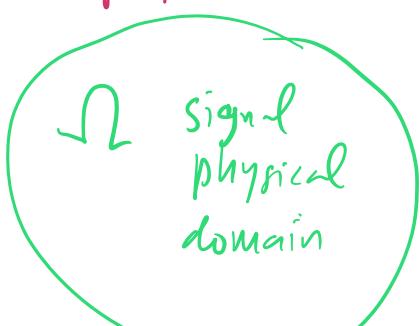
state of N particles

particle physics:

$$\Omega = \{(p_i, q_i) \mid i = 1 \dots N\}$$

$$= (\mathbb{R}^3 \times \mathbb{R}^3)^N$$

→ Q: How to exploit our prior knowledge of the physical domain?



low-dimensional

high-dim

Very-high-dim.

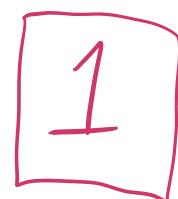
→ Symmetries are known transformations $T: X \rightarrow X$ that do not change "much" the function of interest f :

$T: X \rightarrow X$ such that $f(Tx) \approx f(x)$ $\forall x \in X$

↳ symmetries might be exact or approximate.

$$f(Tx) = f(x)$$

$(f(Tx) = f(x) \text{ yet } \|Tx - x\| = O(\|x\|))$

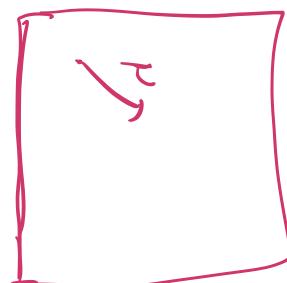


→ Rather than describing symmetries in the data space (high-dim), we rely on the physical space Ω instead.

→ Example: Let $\Omega = [0, 1]^2$ and $\tau \in \Omega$.

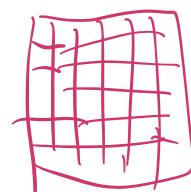
$$T_\tau: \Omega \rightarrow \Omega$$

$$u \mapsto u - \tau \pmod{1}$$



More generally, in a discrete domain $\Omega = [1, n]^2$, we can consider transformations that move pixels around:

$$T: [1, n]^2 \rightarrow [1, n]^2$$



if $|\Omega| = N$, T can be identified with a

permutation $\sigma \in S = \{\pi: \mathbb{N} \rightarrow \mathbb{N}\}$ bijection.

→ A transformation of the domain \mathcal{U} can be "lifted" to a linear transformation on X :

$$\begin{array}{ccc} T: \mathcal{U} \rightarrow \mathcal{U} & \rightsquigarrow & T: X \rightarrow X \\ u \mapsto Tu & & x \mapsto (Tx): \mathcal{U} \rightarrow \mathbb{C} \end{array}$$

Example: $T_\tau: \mathcal{U} \rightarrow \mathcal{U}$

$$u \mapsto u + \tau$$

$$x \mapsto (Tx)(u) = x(u - \tau)$$

$$\left\{ \begin{array}{l} u \mapsto u + \tau \\ \tilde{T}_\tau: X \rightarrow X \end{array} \right.$$

$$u \mapsto (Tx)(u) = x(T^{-1}u)$$

→ observe that \tilde{T} is always linear.

→ The set of exact symmetries of a hypothesis class F (T such that $f(Tx) = f(x) \forall x, \forall f \in F$)

forms a group G with composition: if T and T' are exact symmetries then $T \circ T'$ is also an exact symm.

↳ In this case, the action of the group G on X is called a group representation.

Natural Questions (i) How much can be gained by leveraging symmetries?

(ii) How to implement "symmetry-aware" NNs efficiently? \leftarrow "Geometric Deep Learning"

Sample Complexity of Learning Under Invariance

→ Setup: $\Omega = \{1, \dots, d\}$ discrete domain with d "pixels"

$X = S^{d-1}$ unit d -dimensional sphere; with
 $\nu = \text{Unif}(S^{d-1})$

→ Permutations $\pi : \Omega \rightarrow \Omega$ define unitary linear transformations of $S^{d-1} : \pi$ can be viewed as a $d \times d$ permutation matrix.

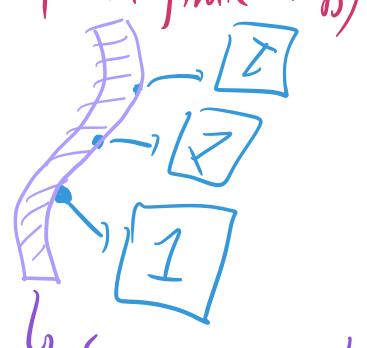
$$\|\pi \cdot x\| = \|x\|.$$

\uparrow
 S^{d-1}

→ Given a subset (or a subgroup, to simplify) G of the permutation group, we define the smoothing operator S_G acting on $L^2(S^{d-1}) = \{f : S^{d-1} \rightarrow \mathbb{R}\}$ (with finite norm)

$$S_G f = \frac{1}{|G|} \sum_{\sigma \in G} f(\sigma \cdot x)$$

(averaging f along each orbit).



\rightarrow f G -invariant function $\left(\begin{array}{l} f \text{ s.t. } f(\sigma x) = f(x) \\ \forall x \in G \end{array} \right) \quad \left\{ \begin{array}{l} \sigma \cdot x ; \sigma \in G \end{array} \right.$

$$S_G f = f$$

\rightarrow f is G -stable if $f = S_G \tilde{f}$ for some $\tilde{f} \in H$

\rightarrow Recall from the spherical harmonic recitation: an orthonormal basis of $L^2(S^{d-1})$

$$\left\{ Y_{kj} \right\}_{k=1 \dots \infty} \leftrightarrow \text{frequency magnitude.}$$

$j=1 \dots N_{d,k} \leftrightarrow \text{frequency "direction".}$

$$N_{d,k} \sim (2k+d-2) \binom{k+d-3}{d-2}$$

$\{Y_{kj}\}_j$ are harmonic polynomials of degree k .

$$f = \sum_{k=1}^{\infty} \sum_{j=1}^{N_{d,k}} \langle f, Y_{j,k} \rangle \cdot Y_{j,k}$$

\hookrightarrow non-linear function.

\rightarrow Dot-product kernels $K(x, x') = \varphi(\langle x, x' \rangle)$ define an RKHS $H_\varphi \subset L^2(S^{d-1})$, given by functions

$$H_\varphi = \left\{ f = \sum_{j,k} \alpha_{j,k} Y_{j,k} \mid \sum_{k,j} \mu_k^{-1} \alpha_{j,k}^2 < +\infty \right\}$$

$$\mu_k = C_d \cdot \int_0^1 \varphi(t) \cdot \left(P_{k,d}(t) \sqrt{(1-t^2)} \right)^{(d-3)} dt$$

\hookrightarrow Legendre / Gegenbauer polynomials.

$$P_{k,d}: [-1, 1] \rightarrow \mathbb{R}$$

$$\ell: [-1, 1] \rightarrow \mathbb{R}.$$

Fact: [Mei, Mikrasenich, Montanari'21] The smoothing operator leaves each harmonic space $V_{lc} = \text{span}\{Y_{l,j} h_j\}$ invariant: $S_G V_{lc} = \tilde{V}_{lc} \subseteq V_K$.

\tilde{V}_{lc} : harmonic polynomials of degree lc that are G -invariant.

$\bar{N}_{d,lc}$: dimension of \tilde{V}_{lc} .

Fact: [MMPM'21]: We have

$$\bar{N}_{d,lc} = N_{d,lc} \left[\frac{1}{|G|} \sum_{\Gamma \in G} E_{x \sim \Gamma} \left[P_{d,lc}(\langle x, \Gamma \cdot x \rangle) \right] \right]$$

$$\textcircled{1} + \sum_{\Gamma \neq Id}$$

$\gamma_{d,lc} :=$ fraction of degree- lc harmonic polynomials that are G -invariant.

$$(\text{Moral: } \gamma_{d,lc} \sim \frac{1}{|G|})$$

Q: If our goal is to learn f' , and we have the promise that f' is G -invariant, how can

we adapt the kernel?

→ G -invariant kernel associated with K_φ :

$$K_G^\varphi(x, x') = \frac{1}{|G|} \sum_{g \in G} \varphi(\langle x, gx' \rangle)$$

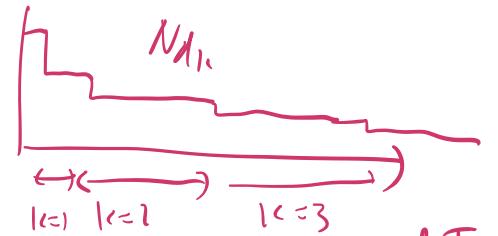
Remark: This is not a user-friendly formula!

→ Associated Integral operators:

$$T_{1c} f(x) = \int K(x, x') f(x') \nu(dx')$$

$$T_{1G} f(x) = \int K_G(x, x') f(x') \nu(dx').$$

spectrum of T_{1c}



spectrum of T_{1G}



↳ T_{1c} and T_{1G} share the same eigenfunctions/eigenvalues;
only difference is the multiplicity
($N_{d1c} \rightarrow \overline{N_{d1c}}$).

→ Let's assume that f^* is G -invariant. Then
the approximation error is not affected by smoothing:

$$\inf \quad \|f - f^*\|^2 + \lambda \|f\|_2^2$$

11

$$\inf_{f \in H_{KG}} \|f - f^*\|^2 + \lambda \|f\|_{H_{KG}}^2$$

→ Generalisation error of kernel ridge regression can be studied using the degrees of freedom:

$$N_K(\lambda) = \text{Tr} \left(\sum_K (\sum_K + \lambda I)^{-1} \right) = \sum_m \frac{\lambda_m}{\lambda_m + \lambda}$$

$\{\lambda_m\}_n$ are the eigenvalues of T_{Kc} .

$N_K(\lambda)$ measures "how big" is the RICS.

→ Punchline: Invariant Kernel has smaller degrees of freedom

$$N_{KG}(\lambda) \lesssim \underbrace{\left(\sup_K \delta_d(k) \right)}_{\hookrightarrow \text{fraction of invariant harmonic polynomials}} \cdot N_K(\lambda)$$

→ Theorem [Bietti, B, Venturi'21]

Assume (i) $f^* = T_K^r g$ for $\|g\|_{L^2} \leq C$ (source cond)

(ii) $N_K(\lambda) \leq C_K \lambda^{-\alpha}$ $\alpha > 1$ (capacity cond)

(iii) f^* is κ -invariant. Then

$$E \| \hat{f}_{\lambda, G} - f^* \| \asymp \left(\frac{1}{|G| \cdot n} \right)^{\frac{2\alpha r}{2\alpha r + 1}} \text{ using } K_G,$$

instead of $E \| \hat{f}_\lambda - f^* \| \asymp \left(\frac{1}{n} \right)^{\frac{2\alpha r}{2\alpha r + 1}}$

Remarks: \oplus Confirm that learning with K_α is more sample efficient. $n \rightsquigarrow n_{\text{eff}} = n \cdot |G|$

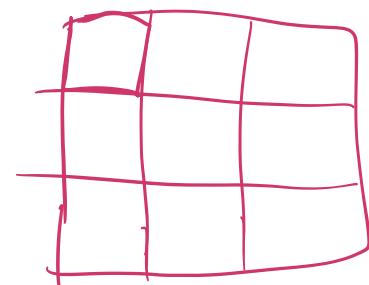
$|G|$ can be either

small ($|G|=d$ for 1d-shifts)

or potentially exponential

in dimension

(local translations).



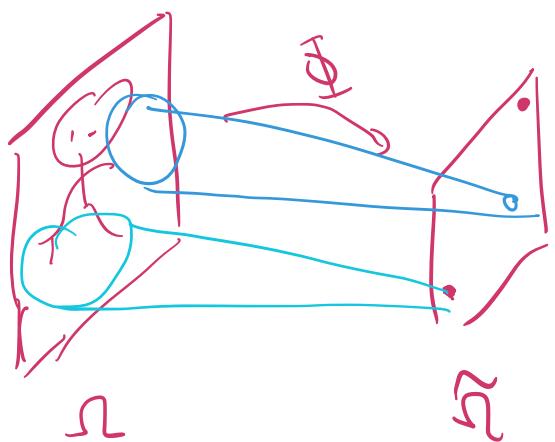
\ominus Invariance cannot change the rate. In particular, if f^* is Lipschitz and invariant,

$$\frac{2\alpha r}{2\alpha r + 1} \sim \frac{1}{d}$$

$\left(\frac{1}{|G|n} \right)^{1/d}$, so learning is still cursed by dimension, even when $|G|$ is $\exp(d)$.

\rightarrow Conclusion: To break the curse of dim, we need another prior on top of symmetry/invariance.

→ What is missing? Scale separation / Multiscale prior.



"composability"

$$f \approx \tilde{f} \circ \phi$$
$$\approx$$