

Geometry of Linear Neural Networks: Equivariance and Invariance under Permutation Groups

Kathlén Kohn, *Anna-Laura Sattelberger*, Vahid Shahverdi
arxiv:2309.13736 [cs.LG]

DMV Annual Meeting 2023
Section: Mathematics of Data Science

September 26, 2023
TU Ilmenau

Two questions

- ① How do the network's properties affect the geometry of its **function space**?
How to characterize **equivariance** or **invariance**?
- ② How to **parameterize** equivariant and invariant networks?
Which implications does it have for **network design**?

Neural networks

A neural network F of depth L is a **parameterized family of functions** $(f_{L,\theta}, \dots, f_{1,\theta})$

$$F: \mathbb{R}^N \longrightarrow \mathcal{F}, \quad F(\theta) = f_{L,\theta} \circ \dots \circ f_{1,\theta} =: f_\theta.$$

Each layer $f_{k,\theta}: \mathbb{R}^{d_{k-1}} \rightarrow \mathbb{R}^{d_k}$ is a composition **activation** \circ **affine-linear**.

Training a network

Given training data $\mathcal{D} = \{(\hat{x}_i, \hat{y}_i)_{i=1, \dots, S}\} \subset \mathbb{R}^{d_0} \times \mathbb{R}^{d_L}$, the aim is to minimize the loss

$$\mathcal{L}: \mathbb{R}^N \xrightarrow{F} \mathcal{F} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}.$$

Example: For $\ell_{\mathcal{D}}$ the squared error loss, this gives $\min_{\theta \in \mathbb{R}^N} \sum_{i=1}^S (f_\theta(\hat{x}_i) - \hat{y}_i)^2$.

On function space: $\min_{M \in \mathcal{F}} \|M\hat{X} - \hat{Y}\|_{\text{Frob}}^2$.

Critical points of \mathcal{L}

- ◇ **pure:** critical point of $\ell_{\mathcal{D}}$
- ◇ **spurious:** induced by parameterization

Linear convolutional networks (LCNs)

- ◇ **linear**: identity as activation function
- ◇ **convolutional** layers with filter $w \in \mathbb{R}^k$ and stride $s \in \mathbb{N}$:

$$\alpha_{w,s}: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}, \quad (\alpha_{w,s}(x))_i = \sum_{j=0}^{k-1} w_j x_{is+j}.$$

Geometry of linear convolutional networks [1]

Function space $\mathcal{F}_{(d,s)}$ of LCN: semi-algebraic set, Euclidean-closed

Theorem [2]

Let (\mathbf{d}, \mathbf{s}) be an LCN architecture with all strides > 1 and $N \geq 1 + \sum_i d_i s_i$. For almost all data $\mathcal{D} \in (\mathbb{R}^{d_0} \times \mathbb{R}^{d_L})^N$, every critical point θ_c of \mathcal{L} satisfies one of the following:

- 1 $F(\theta_c) = 0$, or
- 2 θ_c is a regular point of F and $F(\theta_c)$ is a **smooth, interior point** of $\mathcal{F}_{(d,s)}$.
In particular, $F(\theta_c)$ is a critical point of $\ell_{\mathcal{D}}|_{\text{Reg}(\mathcal{F}_{(d,s)})}$.

This is known to be false for . . .

- ◇ linear fully-connected networks
- ◇ stride-one LCNs

[1] K. Kohn, T. Merkh, G. Montúfar, M. Trager. Geometry of Linear Convolutional Networks. *SIAM J. Appl. Algebra Geom.*, 6(3):368–406, 2022.

[2] K. Kohn, G. Montúfar, V. Shahverdi, M. Trager. Function Space and Critical Points of Linear Convolutional Networks. Preprint arXiv:2304.0572, 2023.

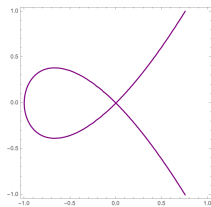
Natural points of entry

- ◇ algebraic vision [3]
- ◇ geometry of function spaces

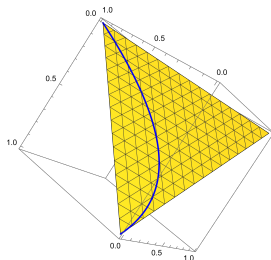
Algebraic varieties

subsets of \mathbb{C}^n obtained as common **zero set of polynomials** $p_1, \dots, p_N \in \mathbb{C}[x_1, \dots, x_n]$

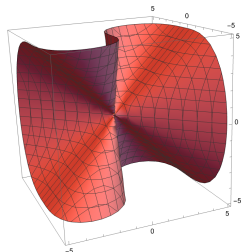
Drawing real points of algebraic varieties



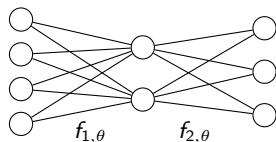
$\mathcal{V}(y^2 - x^2(x + 1))$
a nodal curve



$\mathcal{V}(p_0 p_2 - (p_0 + p_1) p_1) \subset \Delta_2$
a discrete statistical model



$\mathcal{V}(x^2 y - y^3 - z^3)$
a cubic surface



Example

$$F: \mathbb{R}^{2 \times 4} \times \mathbb{R}^{3 \times 2} \longrightarrow \mathbb{R}^{3 \times 4}, \quad (M_1, M_2) \mapsto M_2 \cdot M_1$$

$$\text{parameter space: } \mathbb{R}^N = \mathbb{R}^{2 \times 4} \times \mathbb{R}^{3 \times 2}, \quad f_{1,\theta} = M_1, \quad f_{2,\theta} = M_2$$

Its function space \mathcal{F} is the set of real points of the **determinantal variety**

$$\mathcal{M}_{2,3 \times 4}(\mathbb{R}) = \left\{ M \in \mathbb{R}^{3 \times 4} \mid \text{rank}(M) \leq 2 \right\}.$$

The determinantal variety $\mathcal{M}_{r,m \times n}$

For $M = (m_{ij})_{i,j} \in \mathbb{C}^{m \times n}$: $\text{rank}(M) \leq r \Leftrightarrow$ all $(r+1) \times (r+1)$ minors of M vanish.
 Define **polynomials in m_{ij}**

$$\mathcal{M}_{r,m \times n} = \{ M \mid \text{rank}(M) \leq r \} \subset \mathbb{C}^{m \times n}.$$

Well studied! $\dim(\mathcal{M}_{r,m \times n}) = r(m+n-r)$, $\mathcal{M}_{r,m \times n}(\mathbb{R})$, singularities, ...

$f_\theta: \mathbb{R}^n \rightarrow \mathbb{R}^r \rightarrow \mathbb{R}^m$ $r < \min(m, n)$
 $G = \langle \sigma_1, \dots, \sigma_g \rangle \leq \mathcal{S}_n$ a permutation group, acting on \mathbb{R}^n by permuting the entries
induced action on M : permuting its columns

Invariance under $\sigma \in \mathcal{S}_n$: $f_\theta \circ \sigma \equiv f_\theta$

Decomposing into cycles

The decomposition $\sigma = \pi_1 \circ \dots \circ \pi_k$ of σ into k disjoint cycles induces a partition $\mathcal{P}(\sigma) = \{A_1, \dots, A_k\}$ of the set $[n] = \{1, \dots, n\}$. $A_1, \dots, A_k \subset [n]$ pairwise disjoint sets

Example: The permutation $\sigma = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 3 & 5 & 4 & 1 & 2 \end{pmatrix} = (134)(25) \in \mathcal{S}_5$ induces the partition $\mathcal{P}(\sigma) = \{\{1, 3, 4\}, \{2, 5\}\}$ of $[5] = \{1, 2, 3, 4, 5\}$. For $\eta = (143)(25) \neq \sigma$: $\mathcal{P}(\eta) = \mathcal{P}(\sigma)$.

Characterizing invariance $MP_\sigma \stackrel{!}{=} M$

Let $\sigma \in \mathcal{S}_n$ and $\mathcal{P}(\sigma) = \{A_1, \dots, A_k\}$ its induced partition. A matrix $M = (m_1 | \dots | m_n)$ is invariant under $\sigma = \pi_1 \circ \dots \circ \pi_k$ if and only if for each i , the columns $\{m_j\}_{j \in A_i}$ coincide.

\Rightarrow If M is invariant under σ , its rank is at most k .

Properties of $\mathcal{I}_{r,m \times n}^G \subset \mathcal{M}_{r,m \times n}$

$G = \langle \sigma_1, \dots, \sigma_g \rangle \leq \mathcal{S}_n$ a permutation group
 $\sigma_i = \pi_{i,1} \circ \dots \circ \pi_{i,k_i}, i = 1, \dots, g$ decomposition into pairwise disjoint cycles π_i

Reduction to cyclic case

There exists $\sigma \in \mathcal{S}_n$ such that $\mathcal{I}_{r,m \times n}^G = \mathcal{I}_{r,m \times n}^\sigma$. Any σ for which $\mathcal{P}(\sigma)$ is the **finest common coarsening** of $\mathcal{P}(\sigma_1), \dots, \mathcal{P}(\sigma_g)$ does the job!

Proposition

Let $G = \langle \sigma \rangle \leq \mathcal{S}_n$ be cyclic, and $\sigma = \pi_1 \circ \dots \circ \pi_k$ its decomposition into pairwise disjoint cycles π_i . The variety $\mathcal{I}_{r,m \times n}^\sigma$ is isomorphic to the determinantal variety $\mathcal{M}_{\min(r,k), m \times k}$ via a linear isomorphism $\psi_{\mathcal{P}(\sigma)}: \mathcal{I}_{r,m \times n}^\sigma \rightarrow \mathcal{M}_{\min(r,k), m \times k}$. deleting repeated columns

Via that, we can determine $\dim(\mathcal{I}_{r,m \times n}^\sigma)$, $\deg(\mathcal{I}_{r,m \times n}^\sigma)$, and $\text{Sing}(\mathcal{I}_{r,m \times n}^\sigma)$.

Example ($m = 2, n = 5, r = 1$)

Let $\sigma = (134)(25) \in \mathcal{S}_5$ and hence $k = 2$. Any invariant matrix $M \in \mathcal{M}_{2 \times 5}(\mathbb{R})$ is of the form $\begin{pmatrix} a & c & a & a & c \\ b & d & b & b & d \end{pmatrix}$ for some $a, b, c, d \in \mathbb{R}$. The rank constraint $r = 1$ imposes that $(c, d) = \lambda \cdot (a, b)^\top$ for some $\lambda \in \mathbb{R}$, where we assume that $(a, b) \neq (0, 0)$. Then

$$\psi_{\mathcal{P}(\sigma)}: \begin{pmatrix} a & \lambda a & a & a & \lambda a \\ b & \lambda b & b & b & \lambda b \end{pmatrix} \mapsto \begin{pmatrix} a & \lambda a \\ b & \lambda b \end{pmatrix}.$$

Parameterizing invariance and network design

$\mathcal{S}_n \ni \sigma = \pi_1 \circ \dots \circ \pi_k$, $\mathcal{P}(\sigma) = \{A_1, \dots, A_k\}$

Invariance of $M \in \mathcal{M}_{m \times n}$: forces columns $\{m_j\}_{j \in A_i}$ to coincide. For each i , remember representative m_{A_i} and denote $M_1 := (m_{A_1} \mid \dots \mid m_{A_k}) \in \mathcal{M}_{m \times k}$.

Parameterization

Any σ -invariant $M \in \mathcal{M}_{m \times n}$ of rank k factorizes as $M = M_1 \cdot (e_{i_1} \mid \dots \mid e_{i_n})$, $e_{i_j} \in \mathbb{R}^k$.
 i -th standard unit vector in column j for all $j \in A_i$

Fibers of multiplication map

Let $r \leq \min(m, n)$. Denote by $p: \mathcal{M}_{m \times r} \times \mathcal{M}_{r \times n}$, $(A, B) \mapsto A \cdot B$. If $\text{rank}(M) = r$ and $M = p(A, B)$ for some A, B , then the fiber of p over M is

$$p^{-1}(M) = \left\{ (AT^{-1}, TB) \mid T \in \text{GL}_r(\mathbb{C}) \right\} \subset \mathcal{M}_{m \times r} \times \mathcal{M}_{r \times n}.$$

Learning invariant linear functions with autoencoders

Let M be invariant under σ and of rank k . **Any** factorization $M = A \cdot B$ is of the form

$$(A, B) \in \left\{ (M_1 T^{-1}, T(e_{i_1} \mid \dots \mid e_{i_n})) \mid T \in \text{GL}_k \right\}.$$

This parameterization imposes a **weight sharing property** on the encoder!

Motivation: complexity during and after training

- 1 For an arbitrary learned function, find a nearest invariant function .
- 2 Training invariant networks: determine pure critical points for Euclidean loss .

Definition

The **Euclidean distance (ED) degree** of an algebraic variety \mathcal{X} in \mathbb{R}^N is the number of complex critical points of the squared Euclidean distance from \mathcal{X} to a general point outside the variety. It is denoted by $\deg_{\text{ED}}(\mathcal{X})$.

Examples: $\deg_{\text{ED}}(\text{circle}) = 2$, $\deg_{\text{ED}}(\text{ellipse}) = 4$.

ED degree of $\mathcal{M}_{r,m \times n}(\mathbb{R})$ and $\mathcal{I}_{r,m \times n}^\sigma(\mathbb{R})$

Let $\sigma = \pi_1 \circ \dots \circ \pi_k \in \mathcal{S}_n$ and $r \leq \min(m, n)$. Then

- ◇ $\deg_{\text{ED}}(\mathcal{M}_{r,m \times n}(\mathbb{R})) = \binom{\min(m,n)}{r}$,
- ◇ $\deg_{\text{ED}}(\mathcal{I}_{r,m \times n}^\sigma(\mathbb{R})) = \deg_{\text{ED}}(\mathcal{M}_{\min(r,k),m \times k}(\mathbb{R})) = \binom{\min(m,k)}{\min(r,k)}$.

Equivariant linear autoencoders

$$f_\theta: \mathbb{R}^n \longrightarrow \mathbb{R}^r \longrightarrow \mathbb{R}^n \quad r < n$$

$G = \langle \sigma \rangle \leq \mathcal{S}_n$ a **cyclic** permutation group generated by a single $\sigma \in \mathcal{S}_n$

Equivariance under σ : $f_\theta \circ \sigma \equiv \sigma \circ f_\theta$.

For matrices: M equivariant if $MP_\sigma = P_\sigma M$. commutator of P_σ

In- and output

- ◇ $n = m^2$: $m \times m$ image with real pixels
- ◇ $n = m^3$: cubic 3D scenery

Characterizing $\mathcal{E}_{r,n \times n}^\sigma$

- ◇ dim: ✓
- ◇ deg: ✓
- ◇ Sing: ✓
- ◇ ED degree: under construction!

Exploiting similarity transforms of the form

$$P_\sigma = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{pmatrix} \xrightarrow{\sim T_1} \left(\begin{array}{ccc|cc} 0 & 0 & 1 & & \\ 1 & 0 & 0 & & 0 \\ 0 & 1 & 0 & & \\ \hline & & & 0 & 1 \\ & & & 1 & 0 \end{array} \right) \xrightarrow{\sim T_2} \begin{pmatrix} 1 & & & & \\ & \zeta_3 & & & \\ & & \zeta_3^2 & & \\ & & & 1 & \\ & & & & -1 \end{pmatrix}.$$

Key points: algebraic geometry helps for...

- 1 a thorough study of function spaces of linear neural networks
fully connected, convolutional
- 2 understanding the training process
locating critical points of the loss
- 3 the design of neural networks
invariance implies rank constraint & weight sharing property
- 4 determining the complexity during and post training
ED degree of real varieties



Future work

- ◇ full characterization of equivariance
non-cyclic permutation groups
- ◇ generalization to other groups
e.g. non-discrete groups
- ◇ variation of the network architecture
more layers, non-linear activation functions

Thank you for your attention!

Characterizing invariance

$\mathcal{S}_n \ni \sigma = \pi_1 \circ \cdots \circ \pi_k$ decomposition of σ into k pairwise disjoint cycles
 $\psi_{\mathcal{P}(\sigma)}: \mathcal{I}_{r,m \times n}^\sigma \cong \mathcal{M}_{\min(r,k), m \times k}$ linear isomorphism

Properties of $\mathcal{I}_{r,m \times n}^\sigma$

$$\begin{aligned} \dim(\mathcal{I}_{r,m \times n}^\sigma) &= \min(r, k) \cdot (m + k - \min(r, k)), \\ \deg(\mathcal{I}_{r,m \times n}^\sigma) &= \prod_{i=0}^{k - \min(r, k) - 1} \frac{(m + i)! \cdot i!}{(\min(r, k) + i)! \cdot (m - (\min(r, k) + i))!}, \\ \text{Sing}(\mathcal{I}_{r,m \times n}^\sigma) &= \psi_{\mathcal{P}(\sigma)}^{-1}(\mathcal{M}_{\min(r, k) - 1, m \times k}). \end{aligned}$$

Euclidean distance degree

$$\deg_{\text{ED}}(\mathcal{I}_{r,m \times n}^G(\mathbb{R})) = \binom{\min(m, k)}{\min(r, k)}.$$

Example

Let $m = n = 5$, $r = 2$ and $\sigma = (134)(25) \in \mathcal{S}_5$. If a matrix $M = AB \in \mathcal{I}_{2,5 \times 5}^\sigma$ is invariant under σ , the encoder factor B has to fulfill the following weight sharing property:

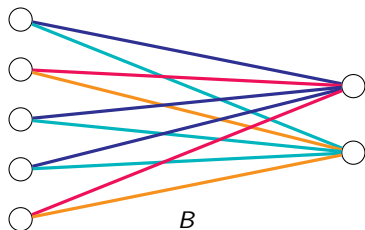


Figure: The σ -weight sharing property imposed on the encoder.

Stepwise diagonalization of permutation matrices: an example

Consider the permutation $\sigma = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 3 & 5 & 4 & 1 & 2 \end{pmatrix} = (1\ 3\ 4)(2\ 5) \in \mathcal{S}_5$. Then

$$P_\sigma = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{pmatrix} \xrightarrow{\sim T_1} \left(\begin{array}{ccc|cc} 0 & 0 & 1 & & \\ 1 & 0 & 0 & 0 & \\ 0 & 1 & 0 & & \\ \hline & & & 0 & 1 \\ & & & 1 & 0 \end{array} \right) \xrightarrow{\sim T_2} \begin{pmatrix} 1 & & & & \\ & \zeta_3 & & & \\ & & \zeta_3^2 & & \\ & & & 1 & \\ & & & & -1 \end{pmatrix}$$

with

$$T_1 = \begin{pmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \quad \text{and} \quad T_2 = \left(\begin{array}{ccc|cc} 1 & 1 & 1 & & \\ 1 & \zeta_3 & \zeta_3^2 & 0 & \\ 1 & \zeta_3^2 & \zeta_3 & & \\ \hline & & & 1 & 1 \\ & & & 1 & -1 \end{array} \right) \in \text{GL}_5(\mathbb{C}),$$

where ζ_3 denotes the primitive 3rd root of unity $\exp^{2\pi i/3}$.

+ grouping identical eigenvalues (optional step)

N.B.: T_2 is block diagonal with Vandermonde matrix blocks $V(1, \zeta_3, \zeta_3^2)$ and $V(1, -1)$.

For a subvariety $\mathcal{X} \subset \mathcal{M}_{m \times n}$ and any $T \in \mathrm{GL}_n(\mathbb{C})$, we denote by \mathcal{X}^T the image of \mathcal{X} under the linear isomorphism

$$\cdot T: \mathcal{M}_{m \times n} \longrightarrow \mathcal{M}_{m \times n}, \quad M \mapsto MT.$$

Lemma

Let $\mathcal{X} \subset \mathcal{M}_{m \times n}$ be a subvariety and let $T \in \mathrm{GL}_n(\mathbb{C})$. Then, $\dim(\mathcal{X}^T) = \dim \mathcal{X}$, $\deg(\mathcal{X}^T) = \deg \mathcal{X}$, $\mathrm{Sing}(\mathcal{X}^T) = \mathrm{Sing}(\mathcal{X})^T$, and $(\mathcal{X}^T) \cap \mathcal{M}_{r, m \times n} = (\mathcal{X} \cap \mathcal{M}_{r, m \times n})^T$ for any $r \leq \min(m, n)$.

Notation: For $T \in \mathrm{GL}_n(\mathbb{C})$ and $M \in \mathcal{M}_{n \times n}$, denote $M^{\sim T} := T^{-1}MT$.

Observation

A matrix M commutes with a matrix P if and only if $P^{\sim T}$ commutes with $M^{\sim T}$, and $MP = M$ if and only if $M^{\sim T}P^{\sim T} = M^{\sim T}$ if and only if $MTP^{\sim T} = MT$.

Proposition

There is a one-to-one correspondence between the **irreducible components** of $\mathcal{E}_{r,n \times n}^\sigma$ and the integer solution vectors $\mathbf{r} = (r_{l,m})$ of

$$\sum_{l \geq 1} \sum_{m \in (\mathbb{Z}/l\mathbb{Z})^\times} r_{l,m} = r,$$

where $0 \leq r_{l,m} \leq d_l$. d_l the dimension of the eigenspace of P_σ of the eigenvalue $\zeta_k = e^{2\pi i/l}$

Properties of $\mathcal{E}_{r,n \times n}^\sigma$

$$\begin{aligned} \dim(\mathcal{E}_{r,n \times n}^\sigma(\mathbb{C})) &= \max_{\mathbf{r}=(r_{l,m})} \left\{ \sum_{l \geq 1} \sum_{m \in (\mathbb{Z}/l\mathbb{Z})^\times} (2d_k - r_{l,m}) \cdot r_{l,m} \right\}, \\ \deg(\mathcal{E}_{r,n \times n}^{\sigma,(\mathbf{r})}(\mathbb{C})) &= \prod_{l \geq 1} \prod_{m \in (\mathbb{Z}/l\mathbb{Z})^\times} \prod_{i=0}^{d_k - r_{l,m} - 1} \frac{(d_k + i)! \cdot i!}{(r_{l,m} + i)! \cdot (d_k - r_{l,m} + i)!}, \\ \text{Sing}(\mathcal{E}_{r,n \times n}^\sigma(\mathbb{K})) &= \mathcal{E}_{r-1,n \times n}^\sigma(\mathbb{K}). \quad \mathbb{K} \in \{\mathbb{R}, \mathbb{C}\} \end{aligned}$$