# Algebraic Neural Network Theory

**Kathlén Kohn**

KTH
Digital Futures

# deep learning in a nutshell



parameter space

function space

**Goal:**
find $\theta$ that minimize

$\theta$
learnable weights

$f_\theta$
network function

$L(f_\theta)$
loss value

# Example: MLPs ← multilayer perceptrons

$$\alpha_L \circ \sigma \circ \sigma \circ \ldots \circ \sigma \circ \alpha_2 \circ \sigma \circ \sigma \circ \alpha_1$$
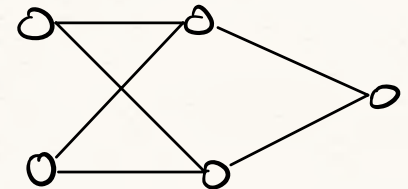
$\alpha_i$ = <u>learnable</u> affine linear functions

$\sigma$ = nonlinear activation function, applied entrywise

we assume: $\sigma$ is a univariate polynomial

**Ex:** $\sigma(x) = x^2$

$$\begin{bmatrix} e & f \end{bmatrix} \sigma \left( \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \right)$$

Which functions does this MLP parametrize?

**Ex:** $\sigma(x) = x^2$

$$\begin{bmatrix} e & f \end{bmatrix} \sigma\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\begin{bmatrix} x \\ y \end{bmatrix}\right)$$



**Which functions does this MLP parametrize?**

$$e(ax+by)^2 + f(cx+dy)^2$$
$$= \underbrace{(a^2e + c^2f)}_{A} x^2 + \underbrace{2(abe + cdf)}_{B} xy + \underbrace{(b^2e + d^2f)}_{C} y^2$$

**Can you obtain _all_ of $R[x,y]_2$ ?**

$\nwarrow$ homogeneous quadratic polynomials in $x,y$

i.e., are all values for $A, B, C$ possible?
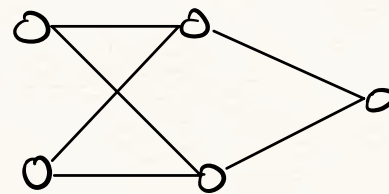
**Ex:** $\sigma(x) = x^2$

$$[e \; f] \; \sigma\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}\right)$$



**Which functions does this MLP parametrize?**

$$e(ax+by)^2 + f(cx+dy)^2$$
$$= \underbrace{(a^2 e + c^2 f)}_{A} x^2 + \underbrace{2(abe + cdf)}_{B} xy + \underbrace{(b^2 e + d^2 f)}_{C} y^2$$

**Can you obtain __all__ of $\mathbb{R}[x,y]_2$ ?**

↖ homogeneous quadratic polynomials in $x,y$

i.e., are all values for $A, B, C$ possible?

YES

**What about $\sigma(x) = x^3$ ?**

**Ex:** $\sigma(x) = x^3$

$$\begin{bmatrix} e & f \end{bmatrix} \sigma\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}\right)$$
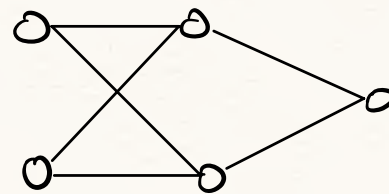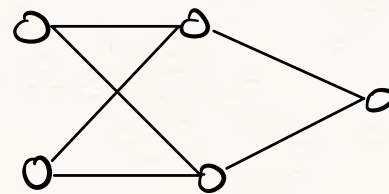


Which functions does this MLP parametrize?

$$e(ax+by)^3 + f(cx+dy)^3$$
$$= \underbrace{(a^3e + c^3f)}_{A}x^3 + \underbrace{3(a^2be + c^2df)}_{B}x^2y + \underbrace{3(ab^2e + cd^2f)}_{C}xy^2 + \underbrace{(b^3e + d^3f)}_{D}y^3$$

Can you obtain __all__ of $\mathbb{R}[x,y]_3$ ?

$\leftarrow$ homogeneous cubic polynomials in $x,y$

i.e., are all values for $A, B, C, D$ possible?

**Ex:** $\sigma(x) = x^3$

$$\begin{bmatrix} e & f \end{bmatrix} \sigma\left( \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \right)$$
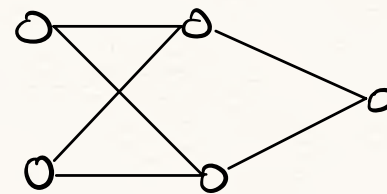


Which functions does this MLP parametrize?

$$e(ax+by)^3 + f(cx+dy)^3$$
$$= \underbrace{(a^3e + c^3f)}_{A} x^3 + \underbrace{3(a^2be + c^2df)}_{B} x^2y + \underbrace{3(ab^2e + cd^2f)}_{C} xy^2 + \underbrace{(b^3e + d^3f)}_{D} y^3$$

Can you obtain **all** of $\mathbb{R}[x,y]_3$ ?

← homogeneous cubic polynomials in $x,y$

i.e., are all values for A, B, C, D possible?

No, e.g.   A = 1
           B = 0
           C = -1
           D = 0

# Neuromanifolds

A parametric machine learning model is a map $\mu: \Theta \times X \longrightarrow Y$.

parameters ↗  inputs ↑  outputs ↖

Its neuromanifold is $\mathcal{M} := \{ \mu(\theta, \cdot): X \to Y \mid \theta \in \Theta \}$.

Example
MLPs:

(no bias vectors
= linear layer functions)



$\sigma(x) = x^2$ $\implies$ $\mathcal{M} = \mathbb{R}[x, y]_2$

$\sigma(x) = x^3$ $\implies$ $\mathcal{M} \subsetneq \mathbb{R}[x, y]_3$

$\sigma(x) = x$ $\implies$ ?

# Neuromanifolds

A parametric machine learning model is a map $\mu: \Theta \times X \to Y$.

parameters ↗   inputs ↑   outputs ↖

Its neuromanifold is $\mathcal{M} := \{ \mu(\theta, \cdot): X \to Y \mid \theta \in \Theta \}$.

Example MLPs:

(no bias vectors = linear layer functions)

$\sigma(x) = x^2$ $\implies$ $\mathcal{M} = \mathbb{R}[x, y]_2$

$\sigma(x) = x^3$ $\implies$ $\mathcal{M} \subsetneq \mathbb{R}[x, y]_3$

$\sigma(x) = x$ $\implies$ $\mathcal{M} = \mathbb{R}^{1 \times 2}$

$\sigma(x) = x$ $\begin{bmatrix} c \\ d \end{bmatrix} \begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$ $\implies$ $\mathcal{M} = ?$
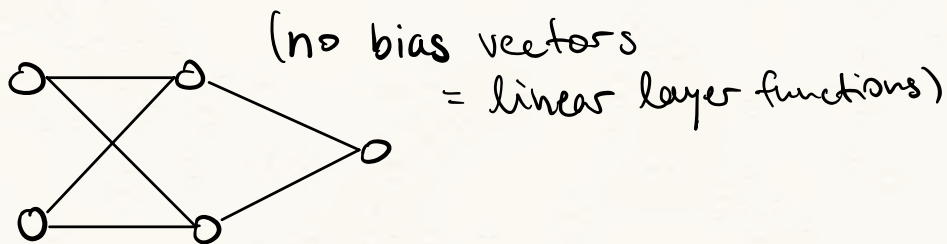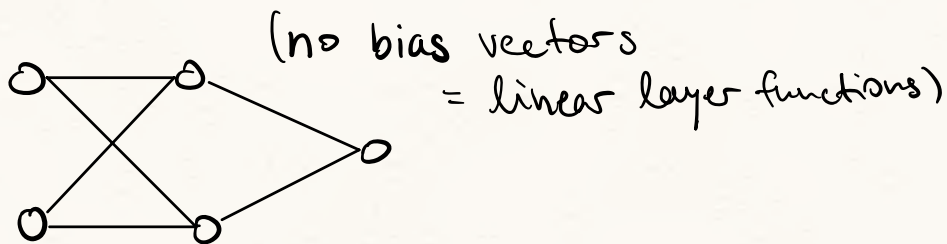
# Neuromanifolds

A **parametric machine learning** model is a map $\mu: \Theta \times X \longrightarrow Y$.

parameters $\nearrow$    inputs $\uparrow$    outputs $\nwarrow$

Its **neuromanifold** is $\mathcal{M} := \{ \mu(\theta, \cdot) : X \to Y \mid \theta \in \Theta \}$.

**Example MLPs:**



(no bias vectors = linear layer functions)
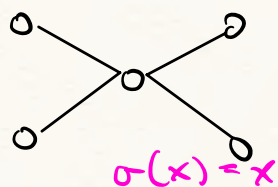
$\sigma(x) = x^2 \qquad \Rightarrow \qquad \mathcal{M} = \mathbb{R}[x, y]_2$

$\sigma(x) = x^3 \qquad \Rightarrow \qquad \mathcal{M} \subsetneq \mathbb{R}[x, y]_3$

$\sigma(x) = x \qquad \Rightarrow \qquad \mathcal{M} = \mathbb{R}^{2 \times 2}$

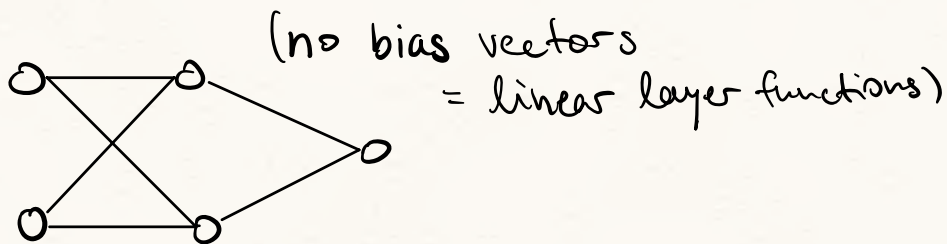$\begin{bmatrix} c \\ d \end{bmatrix} \begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \qquad \Rightarrow \qquad \mathcal{M} = \{ W \in \mathbb{R}^{2 \times 2} \mid \mathrm{rk}(W) \leq 1 \}$

$\sigma(x) = x$

**Linear MLPs:** $\quad \alpha_L \circ \dots \circ \alpha_2 \circ \alpha_1$ , where

$$\alpha_i : \mathbb{R}^{d_{i-1}} \longrightarrow \mathbb{R}^{d_i} \text{ linear}$$

$$\Longrightarrow \mathcal{M} = \{ W \in \mathbb{R}^{d_L \times d_0} \mid rk(W) \leq \min\{d_0, d_1, \dots, d_L\} \}$$

**Linear MLPs:** $\alpha_L \circ \ldots \circ \alpha_2 \circ \alpha_1$, where

$\alpha_i : \mathbb{R}^{d_{i-1}} \to \mathbb{R}^{d_i}$ linear

$\Rightarrow \mathcal{M} = \{ W \in \mathbb{R}^{d_L \times d_0} \mid \mathrm{rk}(W) \leq \min\{d_0, d_1, \ldots, d_L\} \}$

**Polynomial MLPs:** $\alpha_L \circ \sigma \circ \ldots \circ \sigma \circ \alpha_2 \circ \sigma \circ \alpha_1$, where

$\alpha_i : \mathbb{R}^{d_{i-1}} \to \mathbb{R}^{d_i}$ affine linear

$\sigma \in \mathbb{R}[x]_{\leq \delta}$

$\Rightarrow \mathcal{M}$ lives in a $\underline{\text{finite-dimensional}}$ vector space, namely

$$\left( \mathbb{R}[x_1, \ldots, x_{d_0}]_{\leq \delta^{L-1}} \right)^{d_L}$$

**Linear MLPs:** $\quad \alpha_L \circ \dots \circ \alpha_2 \circ \alpha_1$, where

$$\alpha_i : \mathbb{R}^{d_{i-1}} \to \mathbb{R}^{d_i} \text{ linear}$$

$$\Rightarrow \mathcal{M} = \{ W \in \mathbb{R}^{d_L \times d_0} \mid \text{rk}(W) \leq \min\{d_0, d_1, \dots, d_L\} \}$$

**Polynomial MLPs:** $\quad \alpha_L \circ \sigma \circ \dots \circ \sigma \circ \alpha_2 \circ \sigma \circ \alpha_1$, where

$$\alpha_i : \mathbb{R}^{d_{i-1}} \to \mathbb{R}^{d_i} \text{ affine linear}$$

$$\sigma \in \mathbb{R}[x]_{\leq \delta}$$

$\Rightarrow \mathcal{M}$ lives in a <u>finite-dimensional</u> vector space, namely

$$\left( \mathbb{R}[x_1, \dots, x_{d_0}]_{\leq \delta^{L-1}} \right)^{d_L}$$

Polynomial MLPs are the <u>only</u> ones with that property !

# Universal Approximation Theorem

Leshno, Lin, Pinkus, Schocken: Multilayer feedforward networks with a non-polynomial activation function can approximate any function. Neural Networks **6**, 1993 :

## Theorem 1:

Let $\sigma \in M$. Set

$$\Sigma_n = \text{span} \{\sigma(\mathbf{w} \cdot \mathbf{x} + \theta) : \mathbf{w} \in R^n, \theta \in R\}.$$

Then $\Sigma_n$ is dense in $C(R^n)$ if and only if $\sigma$ is not an algebraic polynomial (a.e.).

# Universal Approximation Theorem

Leshno, Lin, Pinkus, Schocken: Multilayer feedforward networks with a non-polynomial activation function can approximate any function. Neural Networks **6**, 1993:
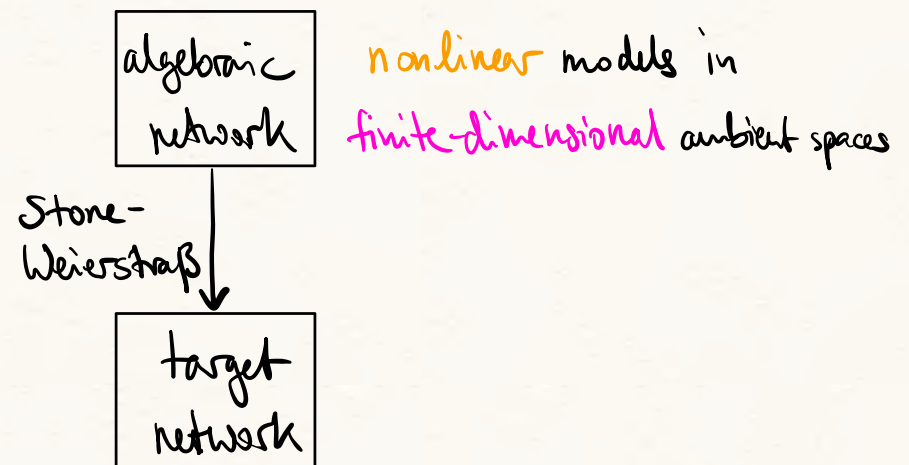
## Theorem 1:

Let $\sigma \in M$. Set

$$\Sigma_n = \text{span} \{\sigma(\mathbf{w} \cdot \mathbf{x} + \theta) : \mathbf{w} \in R^n, \theta \in R\}.$$
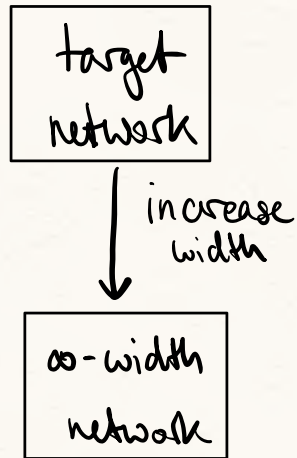
Then $\Sigma_n$ is dense in $C(R^n)$ if and only if $\sigma$ is not an algebraic polynomial (a.e.).

polynomials are <u>the</u> choice to approximate networks with finite-dimensional models
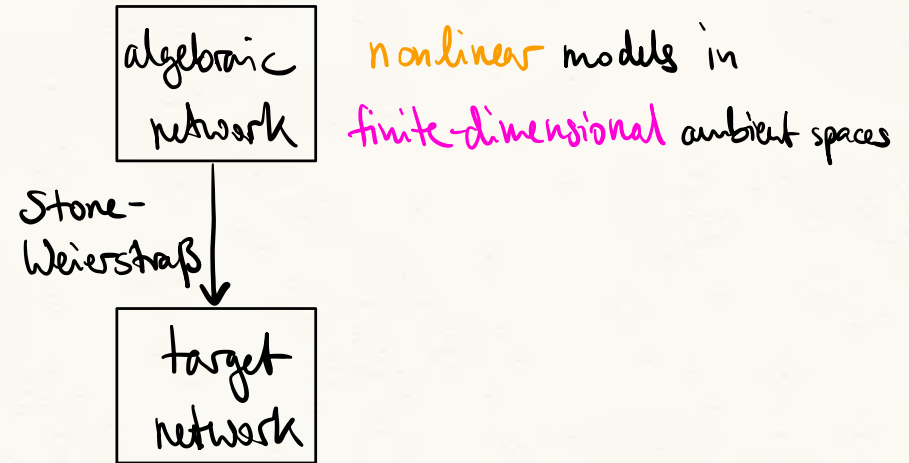
AG approach

nonlinear models in finite-dimensional ambient spaces

| algebraic network |

Stone-Weierstraß

| target network |

neural
tangent
kernel

# NTK approach



target
network

↓ increase
width

∞-width
network

**linearized** models
of **∞ dimension**

algebraic
geometry

# AG approach

algebraic
network

Stone-
Weierstraß ↓

target
network

**nonlinear** models in
**finite-dimensional** ambient spaces

---

**Neural Tangent Kernel: Convergence and Generalization in Neural Networks**

Arthur Jacot, Franck Gabriel, Clement Hongler

Advances in Neural Information Processing Systems 31 (NeurIPS 2018)

> 4000 citations

# Network training = 'distance' minimization

Let $\mathcal{M} \subseteq V := \left( \mathbb{R}[x_1, \ldots, x_{d_0}]_{\leq D} \right)^{d_L}$,

    ↖ neuromanifold

$S \subseteq \mathbb{R}^{d_0} \times \mathbb{R}^{d_L}$ finite dataset,

    ↙ mean squared error

MSE loss: $\mathcal{L}(f) := \sum_{(a,b) \in S} \| f(a) - b \|^2$

    ↙ [dist$(f, g) = 0$ possible for $f \neq g$]

**Proposition:** There is a pseudometric dist$: V \times V \to \mathbb{R}_{\geq 0}$ and some $g \in V$ such that minimizing $\mathcal{L}(f)$ over $f \in \mathcal{M}$ is equivalent to minimizing dist$(f, g)$ over $f \in \mathcal{M}$.

Why?

$V$

# Network training = 'distance' minimization

Let $\mathcal{M} \subseteq V := \left( \mathbb{R}[x_1, \dots, x_{d_0}]_{\leq D} \right)^{d_L}$,
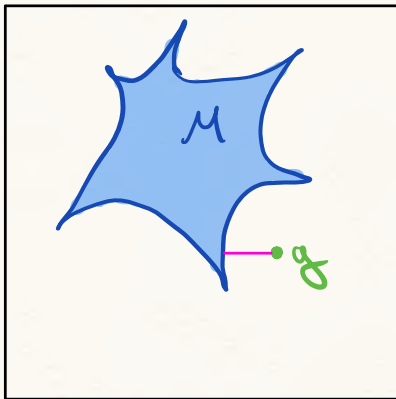
↖ neuromanifold

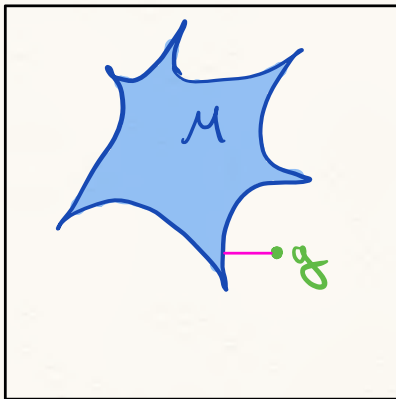$S \subseteq \mathbb{R}^{d_0} \times \mathbb{R}^{d_L}$ finite dataset,

↙ mean squared error

MSE loss: $\mathcal{L}(f) := \sum_{(a,b) \in S} \| f(a) - b \|^2$

↙ [dist$(f, g) = 0$ possible for $f \neq g$]

**Proposition:** There is a pseudometric **dist**$: V \times V \to \mathbb{R}_{\geq 0}$ and some $g \in V$ such that minimizing $\mathcal{L}(f)$ over $f \in \mathcal{M}$ is equivalent to minimizing **dist**$(f, g)$ over $f \in \mathcal{M}$.

$V$



Assume: $d_L = 1$

Veronese embedding ↘

Let $\nu_D: (x_1, \dots, x_{d_0}) \longmapsto$ (all monomials in $x_1, \dots, x_{d_0}$ of degree $\leq D$),

$c_f$ be coefficient vector of $f \in V$ such that $f(x) = \nu_D(x) \cdot c_f$,

# Network training = 'distance' minimization

Let $\mathcal{M} \subseteq V := \left( \mathbb{R}[x_1, \ldots, x_{d_0}]_{\leq D} \right)^{d_L}$,

$\underset{\text{neuromanifold}}{\uparrow}$

$S \subseteq \mathbb{R}^{d_0} \times \mathbb{R}^{d_L}$ finite dataset,

MSE loss: $\underset{\nwarrow \text{mean squared error}}{\quad} \mathcal{L}(f) := \sum_{(a,b) \in S} \| f(a) - b \|^2$

$\underset{\swarrow}{} [\text{dist}(f, g) = 0 \text{ possible for } f \neq g]$

**Proposition:** There is a pseudometric **dist**$: V \times V \to \mathbb{R}_{\geq 0}$ and some $g \in V$ such that minimizing $\mathcal{L}(f)$ over $f \in \mathcal{M}$ is equivalent to minimizing **dist**$(f, g)$ over $f \in \mathcal{M}$.

$V$



Assume: $d_L = 1$

Veronese embedding

Let $\nu_D : (x_1, \ldots, x_{d_0}) \longmapsto$ (all monomials in $x_1, \ldots, x_{d_0}$ of degree $\leq D$), $c_f$ be coefficient vector of $f \in V$ such that $f(x) = \nu_D(x) \cdot c_f$, $A$ & $B$ matrices whose rows are $\nu_D(a)$ & $b$, resp., over all $(a,b) \in S$

$$\Rightarrow \mathcal{L}(f) = \| A c_f - B \|^2$$

# Network training = 'distance' minimization

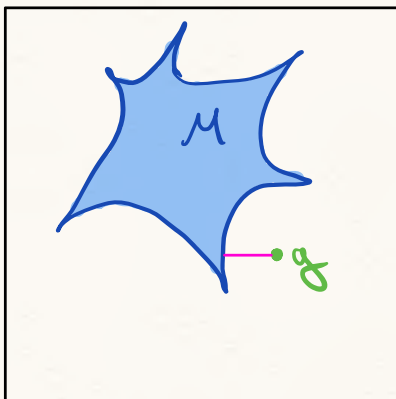Let $\mathcal{M} \subseteq V := \left( \mathbb{R}[x_1, \ldots, x_{d_0}]_{\leq D} \right)^{d_L}$,

$\underset{\text{neuromanifold}}{\uparrow}$

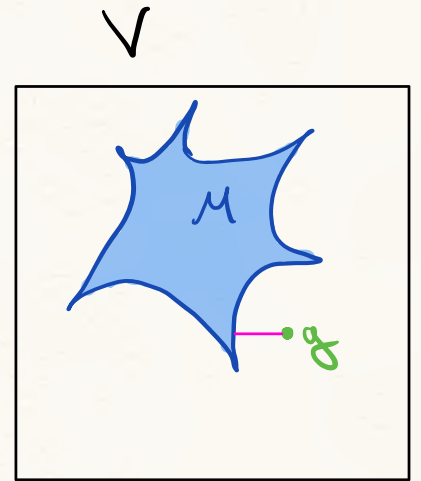$S \subseteq \mathbb{R}^{d_0} \times \mathbb{R}^{d_L}$ finite dataset,

MSE loss: $\quad \mathcal{L}(f) := \sum\limits_{(a,b) \in S} \| f(a) - b \|^2$

$\underset{\text{mean squared error}}{\nwarrow}$

$\left[ \text{dist}(f, g) = 0 \text{ possible for } f \neq g \right]$

**Proposition:** There is a pseudometric **dist** $: V \times V \to \mathbb{R}_{\geq 0}$ and some $g \in V$ such that minimizing $\mathcal{L}(f)$ over $f \in \mathcal{M}$ is equivalent to minimizing **dist**$(f, g)$ over $f \in \mathcal{M}$.

$V$



Assume: $d_L = 1$

Let $\nu_D : (x_1, \ldots, x_{d_0}) \longmapsto$ (all monomials in $x_1, \ldots, x_{d_0}$ of degree $\leq D$),

$\quad$ Veronese embedding $\circlearrowright$

$c_f$ be coefficient vector of $f \in V$ such that $f(x) = \nu_D(x) \cdot c_f$,

$A$ & $B$ matrices whose rows are $\nu_D(a)$ & $b$, resp., over all $(a,b) \in S$

$\Rightarrow \mathcal{L}(f) = \| A c_f - B \|^2 = \| c_f - A^+ B \|^2_{A^T A} + \text{const.}$

$\underset{\text{pseudoinverse}}{\nearrow}$

$\underset{\| c \|_Q := c^T Q c}{\nwarrow}$

$$\underset{f \in M}{\arg\min} \ \mathcal{L}(f) = \underset{f \in M}{\arg\min} \ \| C_f - A^+ B \|^2_{A^T A}$$

$V$



## Observations $(d_L = 1):$

① $A^T A$ depends only on input data,
$A^+ B$ on both input & output

② $A^T A \in \mathbb{R}^{\dim V \times \dim V}$ is rank-deficient whenever $|S| < \dim V \rightsquigarrow \underline{\text{pseudo metric}}$

(LLMs: $|S| < \dim M$)

③

$$\underset{f \in M}{\text{argmin}} \ \mathcal{L}(f) = \underset{f \in M}{\text{argmin}} \ \| C_f - A^+ B \|^2_{A^T A}$$
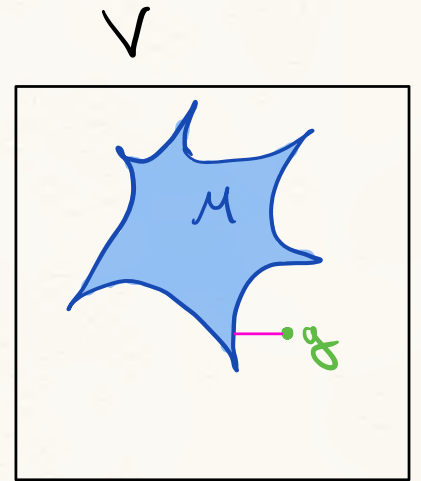
## Observations $(d_L = 1)$:

① $A^T A$ depends only on input data, $A^+ B$ on both input & output

② $A^T A \in \mathbb{R}^{\dim V \times \dim V}$ is rank-deficient whenever $|S| < \dim V$ $\leadsto$ <u>pseudo</u> metric

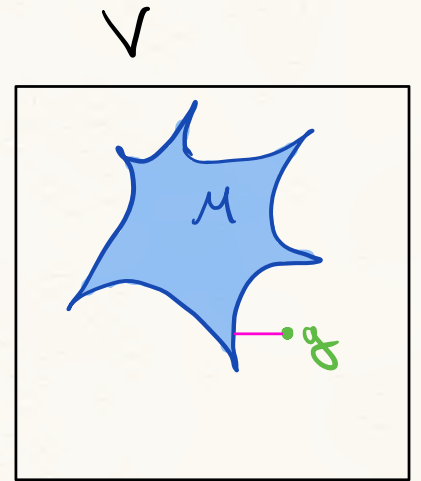(LLMs: $|S| < \dim M$)

③ even when $|S| \gg \dim V$, $A^T A$ is <u>not</u> an arbitrary symmetric PD matrix, while $A^+ B$ yields all vectors $\in \mathbb{R}^{\dim V}$

Why?

Which matrices can be obtained?

(try for $d_0 = 1$: $v(x) = (1, x, x^2, \ldots, x^D)$)

$V$

$$\underset{f \in M}{\arg\min} \; \mathcal{L}(f) = \underset{f \in M}{\arg\min} \; \| C_f - A^+ B \|^2_{A^T A}$$
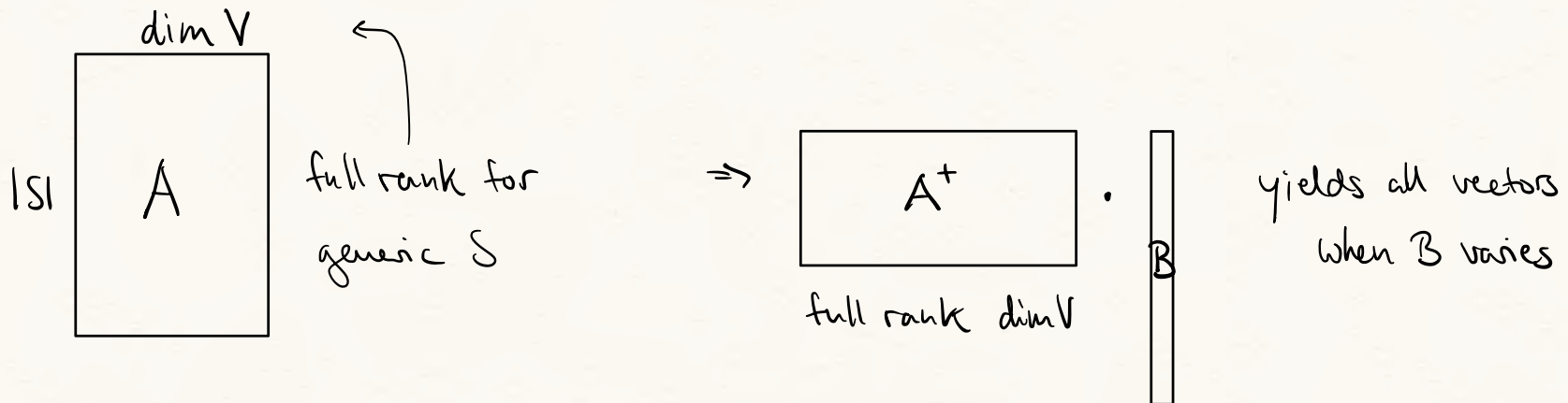

V

## Observations ($d_L = 1$):

① $A^T A$ depends only on input data,
$A^+ B$ on both input & output

② $A^T A \in \mathbb{R}^{\dim V \times \dim V}$ is rank-deficient whenever $|S| < \dim V$ $\leadsto$ <u>pseudo</u> metric

(LLMs: $|S| < \dim M$)

③ even when $|S| \gg \dim V$, $A^T A$ is <u>not</u> an arbitrary symmetric PD matrix,
while $A^+ B$ yields all vectors $\in \mathbb{R}^{\dim V}$



dim V

$|S|$  A    full rank for generic S

$\Rightarrow$   $A^+$ · B   yields all vectors when B varies

full rank dim V

$$\underset{f \in M}{\text{argmin}} \; \mathcal{L}(f) = \underset{f \in M}{\text{argmin}} \; \| C_f - A^+B \|^2_{A^TA}$$
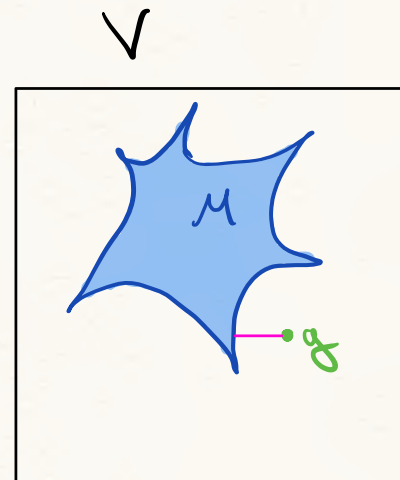

V

## Observations ($d_L = 1$):

① $A^TA$ depends only on input data,
$A^+B$ on both input & output

② $A^TA \in \mathbb{R}^{\dim V \times \dim V}$ is rank-deficient whenever $|S| < \dim V$ ⟿ pseudometric

(LLMs: $|S| < \dim M$)

③ even when $|S| \gg \dim V$, $A^TA$ is **not** an arbitrary symmetric PD matrix,
while $A^+B$ yields all vectors $\in \mathbb{R}^{\dim V}$

$$A^TA = \begin{array}{c} \\ i \rightarrow \end{array} \begin{vmatrix} | & & | \\ v(a_1) & \cdots & v(a_{|S|}) \\ | & & | \end{vmatrix} \begin{vmatrix} \underline{\quad} v(a_1) \underline{\quad} \\ \vdots \\ \underline{\quad} v(a_{|S|}) \underline{\quad} \end{vmatrix}$$

$\downarrow j$

has $(i,j)$ entry $\underset{(a,b) \in S}{\sum} \underbrace{v_i(a) \, v_j(a)}_{\text{monomial of degree } \leq 2D}$

that can be factored in several ways

**Ex.:** $d_0 = 1$

$\Rightarrow v(x) = (1, x, x^2, \ldots, x^D)$

$\Rightarrow A = \begin{bmatrix} 1 & a_1 & a_1^2 & \cdots & a_1^D \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & a_{|S|} & a_{|S|}^2 & \cdots & a_{|S|}^D \end{bmatrix}$    Vandermonde matrix

$\Rightarrow A^T A = \begin{bmatrix} |S| & \sum a_k & \sum a_k^2 & \cdots & \sum a_k^D \\ \sum a_k & \sum a_k^2 & \sum a_k^3 & \cdots & \sum a_k^{D+1} \\ \sum a_k^2 & \sum a_k^3 & \sum a_k^4 & \cdots & \sum a_k^{D+2} \\ \vdots & \vdots & \vdots & & \\ \sum a_k^D & \sum a_k^{D+1} & \sum a_k^{D+2} & \cdots & \sum a_k^{2D} \end{bmatrix}$    Hankel matrix

**Ex.:** $d_0 = 1$

$\Rightarrow v(x) = (1, x, x^2, \ldots, x^D)$

$\Rightarrow A = \begin{bmatrix} 1 & a_1 & a_1^2 & \cdots & a_1^D \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & a_{|S|} & a_{|S|}^2 & \cdots & a_{|S|}^D \end{bmatrix}$  Vandermonde matrix

$\Rightarrow A^T A = \begin{bmatrix} |S| & \sum a_k & \sum a_n^2 & \cdots & \sum a_n^D \\ \sum a_k & \sum a_k^2 & \sum a_k^3 & \cdots & \sum a_k^{D+1} \\ \sum a_k^2 & \sum a_k^3 & \sum a_k^4 & \cdots & \sum a_k^{D+2} \\ \vdots & \vdots & \vdots & & \\ \sum a_k^D & \sum a_k^{D+1} & \sum a_k^{D+2} & \cdots & \sum a_n^{2D} \end{bmatrix}$  Hankel matrix

**Ex.:** $d_0 = 2$, $D = 2$

$\Rightarrow v(x, y) = (1, x, y, x^2, xy, y^2)$

$\Rightarrow A^T A = \sum\limits_{\substack{(a,b) \in S \\ a = (x,y)}} \begin{bmatrix} 1 & x & y & x^2 & xy & y^2 \\ x & x^2 & xy & x^3 & x^2y & xy^2 \\ y & xy & y^2 & x^2y & xy^2 & y^3 \\ x^2 & x^3 & x^2y & x^4 & x^3y & x^2y^2 \\ xy & x^2y & xy^2 & x^3y & x^2y^2 & xy^3 \\ y^2 & xy^2 & y^3 & x^2y^2 & x^3 & y^4 \end{bmatrix} \begin{matrix} 1 \\ x \\ y \\ x^2 \\ xy \\ y^2 \end{matrix}$

# Network training = 'distance' minimization

Let $\mathcal{M} \subseteq V := \left( \mathbb{R}[x_1, \ldots, x_{d_0}]_{\leq D} \right)^{d_L}$,

$\underset{\text{neuromanifold}}{\uparrow}$

$S \subseteq \mathbb{R}^{d_0} \times \mathbb{R}^{d_L}$ finite dataset,

$\underset{\text{mean squared error}}{\nwarrow}$

MSE loss: $\mathcal{L}(f) := \sum_{(a,b) \in S} \| f(a) - b \|^2$

$\underset{\swarrow}{[\text{dist}(f, g) = 0 \text{ possible for } f \neq g]}$

**Proposition:** There is a pseudometric $\text{dist}: V \times V \to \mathbb{R}_{\geq 0}$ and some $g \in V$ such that minimizing $\mathcal{L}(f)$ over $f \in \mathcal{M}$ is equivalent to minimizing $\text{dist}(f, g)$ over $f \in \mathcal{M}$.

$V$



$d_L > 1$

$f = (f_1, \ldots, f_{d_L})$, $\quad C_f := \begin{bmatrix} | & & | \\ c_{f_1} & \cdots & c_{f_{d_L}} \\ | & & | \end{bmatrix}$
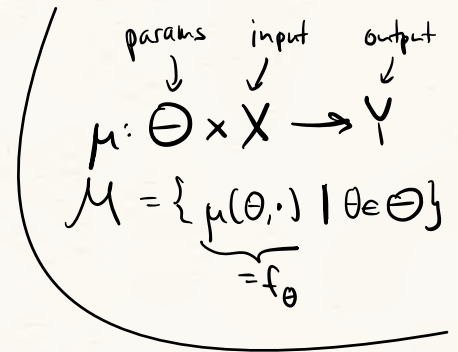
$\Rightarrow f(x) = v_D(x) \cdot C_f$

$\| C \|_Q^2 := \text{tr}(C^T Q C)$

$\Rightarrow \mathcal{L}(f) = \| A C_f - B \|_{\text{Frob}}^2 = \| C_f - A^+ B \|_{A^T A}^2 + \text{const.}$

# Loss Landscape

$$= \{(\theta, \mathcal{L}(f_\theta)) \mid \theta \in \Theta\}$$

params    input    output

$$\mu: \Theta \times X \longrightarrow Y$$

$$\mathcal{M} = \{\underbrace{\mu(\theta, \cdot)}_{= f_\theta} \mid \theta \in \Theta\}$$

# Loss Landscape

$$= \{(\theta, \mathcal{L}(f_\theta)) \mid \theta \in \Theta\}$$

$$
\begin{array}{ccc}
\text{params} & \text{input} & \text{output} \\
\downarrow & \downarrow & \downarrow \\
\mu: \Theta & \times X & \to Y
\end{array}
$$

$$\mathcal{M} = \{\underbrace{\mu(\theta, \cdot)}_{= f_\theta} \mid \theta \in \Theta\}$$

can be studied in a decoupled way:

$$\Theta \xrightarrow{\hspace{3cm}} \mathcal{M} \xrightarrow{\hspace{1cm} \mathcal{L} \hspace{1cm}} \mathbb{R}$$

$$\theta \xmapsto{\hspace{3cm}} f_\theta$$

loss landscape in function space:

$$= \{(f, \mathcal{L}(f)) \mid f \in \mathcal{M}\} \subseteq V \times \mathbb{R}$$

# Loss Landscape

$$= \{(\theta, \mathcal{L}(f_\theta)) \mid \theta \in \Theta\}$$

$$\mu: \underset{\uparrow}{\Theta} \times \underset{\uparrow}{X} \to \underset{\uparrow}{Y}$$
$$\text{params} \quad \text{input} \quad \text{output}$$

$$M = \{\underbrace{\mu(\theta, \cdot)}_{= f_\theta} \mid \theta \in \Theta\}$$

can be studied in a decoupled way:

$$\Theta \longrightarrow M \xrightarrow{\;\;\mathcal{L}\;\;} \mathbb{R}$$
$$\theta \longmapsto f_\theta$$

loss landscape in function space:

$$= \{(f, \mathcal{L}(f)) \mid f \in M\} \subseteq V \times \mathbb{R}$$

Geometry of $M$ affects loss landscape!

How?

Which geometric properties does $M$ have?

# Position: Algebra Unveils Deep Learning
# An Invitation to Neuroalgebraic Geometry

Giovanni Luca Marchetti [*1]  Vahid Shahverdi [*1]  Stefano Mereta [*1]  Matthew Trager [*2]  Kathlén Kohn [*1]

| Machine Learning | Algebraic Geometry |
| --- | --- |
| sample complexity and expressivity | dimension, degree, and covering number |
| subnetworks and implicit bias | singularities |
| identifiability and invariance | fibers of the parameterization |
| optimization and gradient descent | critical point theory, discriminants, and dynamical invariants |

# Identifiability

$$\Theta \xrightarrow{\mu} \mathcal{M}$$

$$\theta \longmapsto f_\theta$$

Given (generic) $f \in \mathcal{M}$, what is $\mu^{-1}(f)$?

# Identifiability

$$\Theta \xrightarrow{\quad \mu \quad} \mathcal{M}$$

$$\theta \longmapsto f_\theta$$

Given (generic) $f \in \mathcal{M}$, what is $\mu^{-1}(f)$?

For **monomial MLP** with $\sigma(x) = x^r$, $r \gg 0$:

$$\mu: \mathbb{R}^{d_L \times d_{L-1}} \times \cdots \times \mathbb{R}^{d_1 \times d_0} \longrightarrow \mathcal{M}$$

$$(W_L, \ldots, W_1) \longmapsto W_L \circ \sigma \circ \cdots \circ \sigma \circ W_2 \circ \sigma \circ W_1$$

What is the generic fiber?

# Identifiability

$$\Theta \xrightarrow{\quad \mu \quad} \mathcal{M}$$

$$\theta \longmapsto f_\theta$$

Given (generic) $f \in \mathcal{M}$, what is $\mu^{-1}(f)$?

For **monomial MLP** with $\sigma(x) = x^r$, $r \gg 0$:

$$\mu : \mathbb{R}^{d_L \times d_{L-1}} \times \cdots \times \mathbb{R}^{d_1 \times d_0} \longrightarrow \mathcal{M}$$

$$(W_L, \ldots, W_1) \longmapsto W_L \circ \sigma \circ \cdots \circ \sigma \circ W_2 \circ \sigma \circ W_1$$

**Observation:** $D_i \in GL(d_i)$ diagonal

$P_i \in GL(d_i)$ permutation matrix

$$\Rightarrow \mu(W_L D_{L-1}^{-r} P_{L-1}^T, \ldots, P_2 D_2 W_2 D_1^{-r} P_1^T, P_1 D_1 W_1)$$

$$= \mu(W_L, \ldots, W_2, W_1)$$

# Identifiability

$$\Theta \xrightarrow{\quad \mu \quad} \mathcal{M}$$
$$\theta \longmapsto f_\theta$$

Given (generic) $f \in \mathcal{M}$, what is $\mu^{-1}(f)$?

For **monomial MLP** with $\sigma(x) = x^r, \ r \gg 0$:

Activation degree thresholds and expressiveness of polynomial neural networks  2024

Bella Finkel,[*] Jose Israel Rodriguez,[†] Chenxi Wu, Thomas Yahl

Proven that those parameter symmetries are the generic fiber (implicitly described all fibers!)

follow-ups (2025):

THE ALEXANDER-HIRSCHOWITZ THEOREM FOR NEUROVARIETIES

ALEX MASSARENTI AND MASSIMILIANO MELLA

Identifiability of Deep Polynomial Neural Networks

Konstantin Usevich,[*] Ricardo Borsoi, Clara Dérand, Marianne Clausel[†]
Université de Lorraine, CNRS, CRAN

# Identifiability

$$\Theta \xrightarrow{\ \mu\ } \mathcal{M}$$

$$\theta \longmapsto f_\theta$$

Given (generic) $f \in \mathcal{M}$, what is $\mu^{-1}(f)$?

For **monomial MLP** with $\sigma(x) = x^r$, $r \gg 0$:

Activation degree thresholds and expressiveness of polynomial neural networks

2024

Bella Finkel,[*] Jose Israel Rodriguez,[†] Chenxi Wu, Thomas Yahl

**Proposition 16.** Let $\mathbb{K}$ be a subfield of $\mathbb{C}$. Given integers $d, k$, there exists an integer $\tilde{r} = \tilde{r}(k)$ with the following property. If $r > \tilde{r}(k)$ and $p_1, \ldots, p_k \in \mathbb{K}[x_1, \ldots, x_d]$ are pairwise non-proportional, then $p_1^r, \ldots, p_k^r$ are linearly independent (over $\mathbb{K}$). Moreover, $\tilde{r}(k) = 6(k-1)^2 - 6(k-1) + 1$ has the desired property.

# Identifiability

$$\Theta \xrightarrow{\quad \mu \quad} \mathcal{M}$$

$$\theta \longmapsto f_\theta$$

Given (generic) $f \in \mathcal{M}$, what is $\mu^{-1}(f)$?

For **polynomial MLP** with $\sigma(x) \in \mathbb{R}[x]_{\leq r}$ generic, $r >> 0$:

**Conjecture:** Generic fiber $\mu^{-1}(f)$ consists only of permutations.

**Conjecture:** Let $d, k \in \mathbb{Z}_{>0}$.

There is $\tilde{r} \in \mathbb{Z}_{>0}$ such that all $r > \tilde{r}$ satisfy:

There is $\mathcal{U} \subseteq \mathbb{R}[x]_{\leq r}$ Zariski open such that:

For all $\sigma \in \mathcal{U}$ and all $p_1, \ldots, p_k \in \mathbb{R}[x_1, \ldots, x_d]$ non-constant & pairwise distinct:

$\sigma(p_1), \ldots, \sigma(p_k)$ are linearly independent.

# Geometry of Neuromanifolds

$$\mu : \Theta \times X \longrightarrow Y \qquad \text{polynomial (in both } \theta \in \Theta \ \& \ x \in X)$$
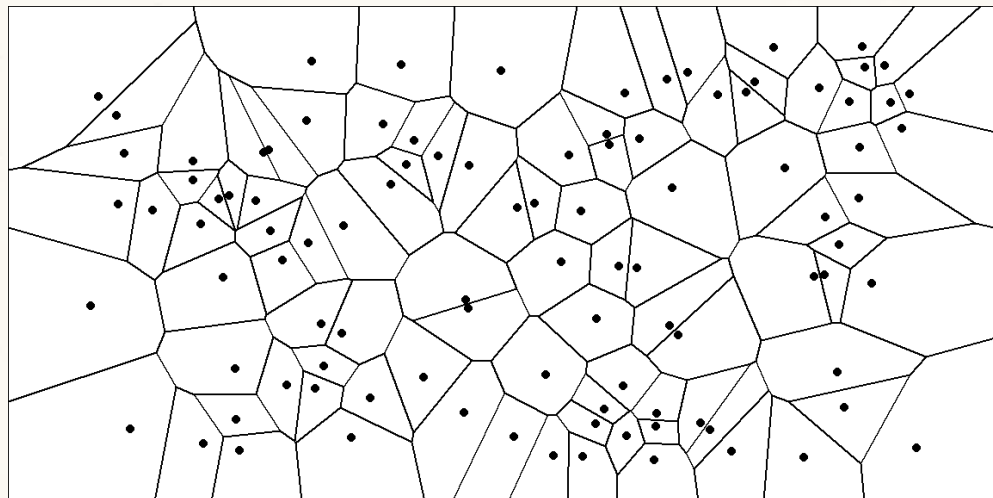
$$\Theta \longrightarrow \mathcal{M}$$
$$\theta \longmapsto \mu(\theta, \cdot)$$

What kind of object is $\mathcal{M}$?

A **semialgebraic** set !
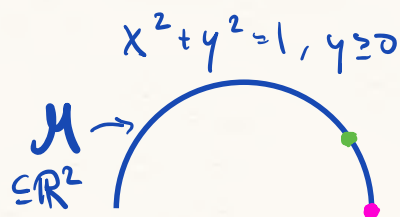
describable by
polynomial equations
& inequalities

# Geometry of Neuromanifolds

$$\mu : \Theta \times X \longrightarrow Y \qquad \text{polynomial (in both } \theta \in \Theta \ \& \ x \in X)$$

$$\Theta \longrightarrow \mathcal{M}$$
$$\theta \longmapsto \mu(\theta, \cdot)$$

What kind of object is $\mathcal{M}$?

A **semialgebraic** set!

↑ describable by polynomial equations & inequalities

Euclidean distance minimization can be implicitly biased to singularities & boundaries of $\mathcal{M}$
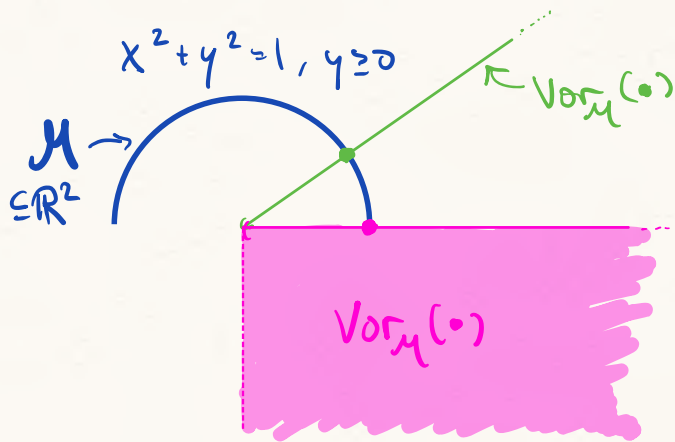
# Voronoi cells



For $S \subseteq \mathbb{R}^n$, the **Voronoi cell** at $p \in S$ is
$$\text{Vor}_S(p) := \{ u \in \mathbb{R}^n \mid \forall q \in S, q \neq p :$$
$$\|p - u\|_2 < \|q - u\|_2 \}$$

$x^2 + y^2 = 1, y \geq 0$

$\mathcal{M} \rightarrow$

$\subseteq \mathbb{R}^2$



What is the Voronoi cell at • ?

What is the Voronoi cell at • ?

# Voronoi cells



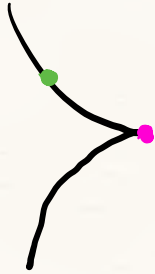For $S \subseteq \mathbb{R}^n$, the **Voronoi cell** at $p \in S$ is
$$Vor_S(p) := \{u \in \mathbb{R}^n \mid \forall q \in S, q \neq p:$$
$$\|p-u\|_2 < \|q-u\|_2\}$$



$x^2 + y^2 = 1, \ y \geq 0$

$\leftarrow Vor_M(\bullet)$

$M \rightarrow$
$\subseteq \mathbb{R}^2$

$Vor_M(\bullet)$

The **2 relative boundary points** are the only points on $M$ with full-dimensional Voronoi cells!
$\rightsquigarrow$ **implicit bias** towards $\partial M$

points in $\partial M$ are global minima with positive probability on data $u$

singularities



What are the Voronoi cells at ● and ● ?

Challenge: Compute this curve!

$$y^2 + x^3 = 0$$

$$t \mapsto (-t^2, t^3)$$

⤳ **implicit bias** towards Sing($M$)

What are the Voronoi cells at • and • ?

singularities

$y^2 + x^3 = 0$

$t \mapsto (-t^2, t^3)$

Challenge: Compute this curve!

⟿ **implicit bias** towards $\text{Sing}(\mathcal{M})$

What are the Voronoi cells at · and · ?

## Tradeoff

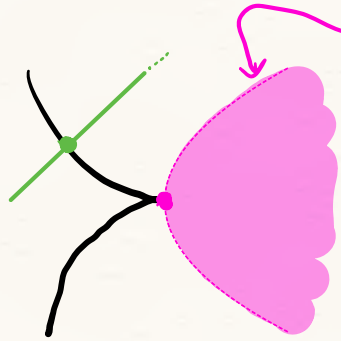learning close to singularity
⟿ slow & numerical instability

[Amari et al]

singular solution generalizes better:
① stable global minimum when perturbing data
② **Conjecture:** singularities of neuromanifolds
are sparse subnetworks
[we've proven this for MLPs & CNNs]

# singularities

$$y^2 + x^3 = 0$$
$$t \longmapsto (-t^2, t^3)$$

**Challenge:** Compute this curve!

⟿ **implicit bias** towards $Sing(\mathcal{M})$

What are the Voronoi cells at • and • ?

## Tradeoff

☹

learning close to singularity
⟿ slow & numerical instability

[Amari et al]

☺

singular solution generalizes better:
① stable global minimum when perturbing data
② **Conjecture:** singularities of neuromanifolds
are sparse subnetworks
[we've proven this for MLPs & CNNs]

In general: depends on **type** of singularity

MLP

CNN

$\sigma(x) =$ generic polynomial of large degree

These singularities have that tradeoff,        .......        while these don't !

In both cases, they are sparse subnetworks ☺

# What about smooth interior points?

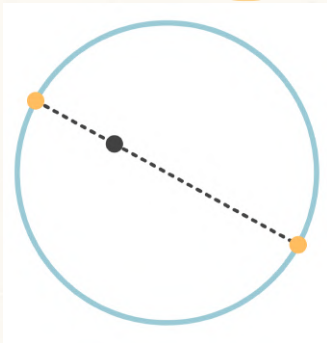$M \subseteq \mathbb{R}^n$ algebraic variety (i.e., described by polynomial equations)

$Q$ symmetric PD $n \times n$ matrix

**Fact:** For almost all $u \in \mathbb{R}^n$, the number of <u>complex</u> critical points of

$$\min_{x \in M \setminus \mathrm{Sing}(M)} \|x - u\|_Q^2$$

is the <u>same</u>, called the **Euclidean Distance Degree**: $EDD_Q(M)$.

What is $EDD_{[\;\cdot\;\cdot\;]}(\bigcirc)$?

# What about smooth interior points?

$M \subseteq \mathbb{R}^n$ algebraic variety (i.e., described by polynomial equations)

$Q$ symmetric PD $n \times n$ matrix

**Fact:** For almost all $u \in \mathbb{R}^n$, the number of <u>complex</u> critical points of

$$\min_{x \in M \setminus Sing(M)} \|x - u\|_Q^2$$

is the <u>same</u>, called the **Euclidean Distance Degree**: $EDD_Q(M)$.

What is $EDD_{\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}}(\bigcirc)$?

What is $EDD_{\begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}}(\bigcirc)$?

# What about smooth interior points?

$M \subseteq \mathbb{R}^n$ algebraic variety (i.e., described by polynomial equations)
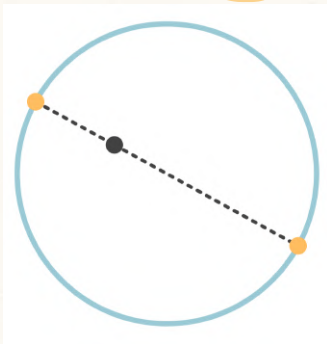
$Q$ symmetric PD $n \times n$ matrix

**Fact:** For almost all $u \in \mathbb{R}^n$, the number of <u>complex</u> critical points of
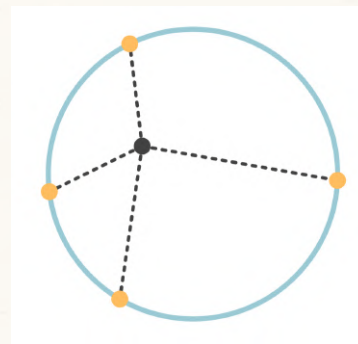
$$\min_{x \in M \setminus \text{Sing}(M)} \|x - u\|_Q^2$$

is the <u>same</u>, called the **Euclidean Distance Degree**: $\text{EDD}_Q(M)$.

What is $\text{EDD}_{\begin{bmatrix}1 & 0 \\ 0 & 1\end{bmatrix}}(\bigcirc)$?



What is $\text{EDD}_{\begin{bmatrix}4 & 0 \\ 0 & 1\end{bmatrix}}(\bigcirc)$?

# What about smooth interior points?

$M \subseteq \mathbb{R}^n$ algebraic variety (i.e., described by polynomial equations)
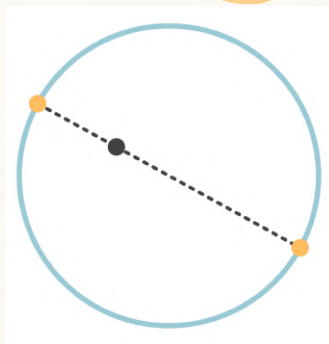
$Q$ symmetric PD $n \times n$ matrix

**Fact:** For almost all $u \in \mathbb{R}^n$, the number of <u>complex</u> critical points of
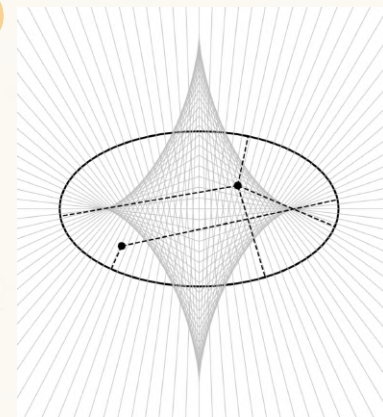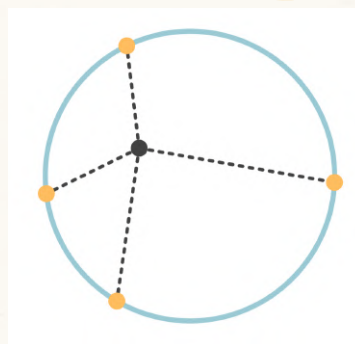
$$\min_{x \in M \setminus \text{Sing}(M)} \|x - u\|_Q^2$$

is the <u>same</u>, called the **Euclidean Distance Degree**: $\text{EDD}_Q(M)$.

What is $\text{EDD}_{\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}}(\bigcirc)$?

What is $\text{EDD}_{\begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}}(\bigcirc)$?
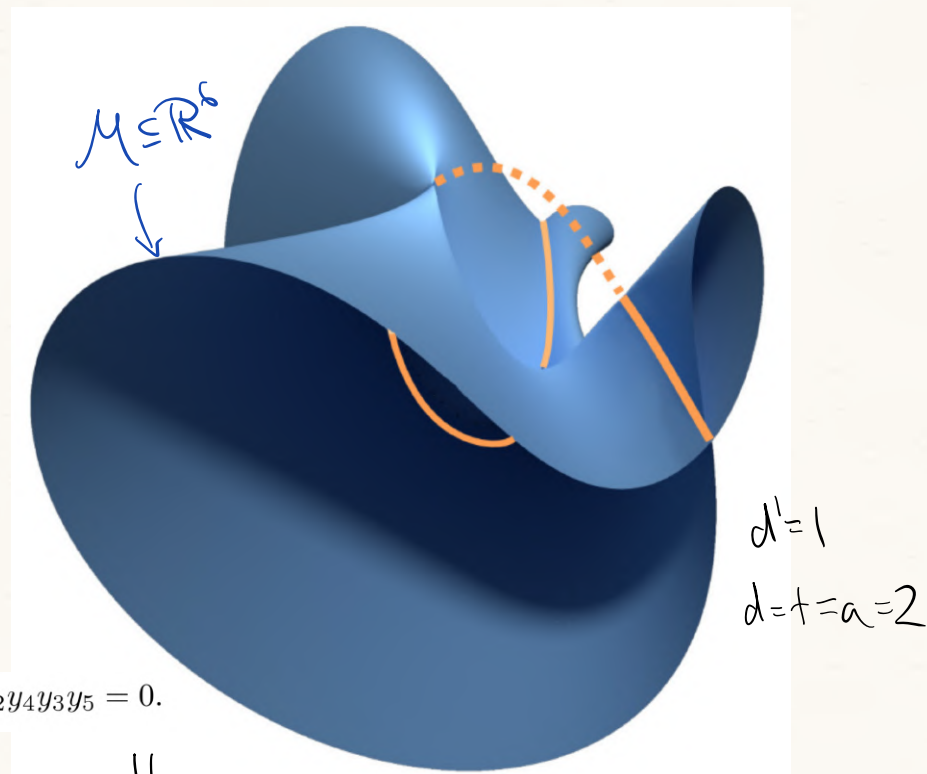
# Lightning Self-Attention (single head, single layer)

$$\mathbb{R}^{d\times t} \longrightarrow \mathbb{R}^{d'\times t}$$
$$X \longmapsto V X X^{\mathsf{T}} K^{\mathsf{T}} Q X$$

learnable parameters
$V \in \mathbb{R}^{d'\times d}$, $K, Q \in \mathbb{R}^{a\times d}$



$M \subseteq \mathbb{R}^6$

$d' = 1$
$d = t = a = 2$

$$y_1^2 y_6^2 + y_4^2 y_3^2 + y_1 y_3 y_5^2 + y_2^2 y_4 y_6 - 2 y_1 y_4 y_3 y_6 - y_2 y_1 y_6 y_5 - y_2 y_4 y_3 y_5 = 0.$$

For almost all PD matrices $Q$,
$$\mathrm{EDD}_Q(M) = 14.$$

What happens if $Q$ becomes degenerate?
(i.e., $Q$ is symmetric positive semidefinite)

# Lightning Self-Attention (single head, single layer)

$$\mathbb{R}^{d \times t} \longrightarrow \mathbb{R}^{d' \times t}$$

$$X \longmapsto V X X^T K^T Q X$$

learnable parameters
$V \in \mathbb{R}^{d' \times d}$, $K, Q \in \mathbb{R}^{a \times d}$

$M \in \mathbb{R}^6$

$d' = 1$
$d = t = a = 2$

$$y_1^2 y_6^2 + y_4^2 y_3^2 + y_1 y_3 y_5^2 + y_2^2 y_4 y_6 - 2 y_1 y_4 y_3 y_6 - y_2 y_1 y_6 y_5 - y_2 y_4 y_3 y_5 = 0.$$

For almost all PD matrices $Q$,
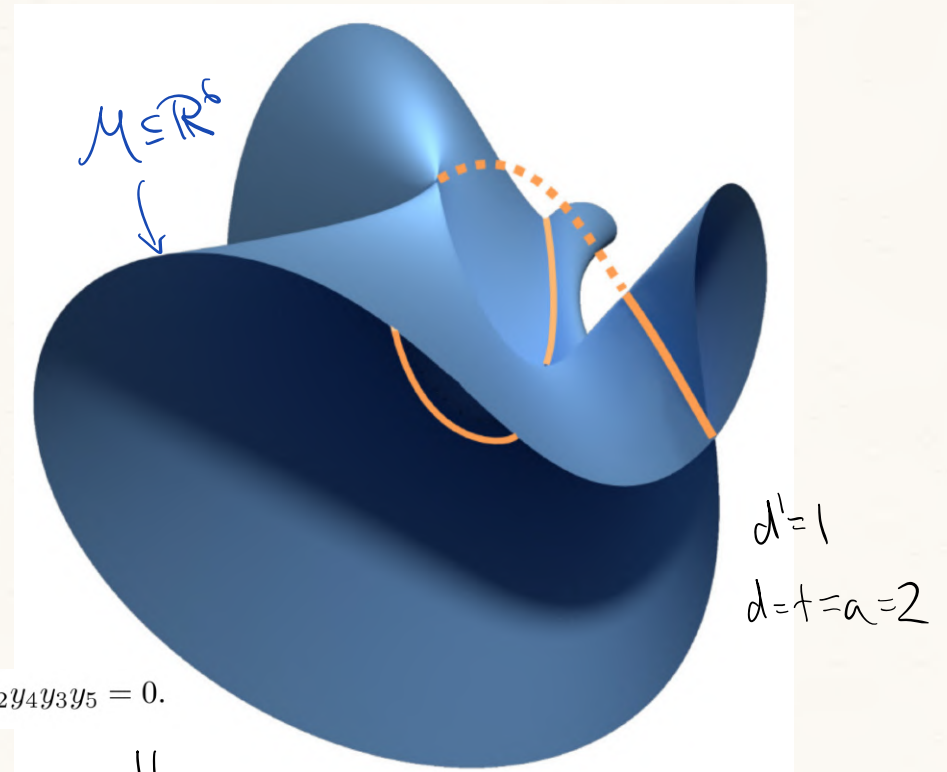$$EDD_Q(M) = 14.$$

What happens if $Q$ becomes degenerate?
(ie., $Q$ is symmetric positive semidefinite)

$K := \dim \ker Q$

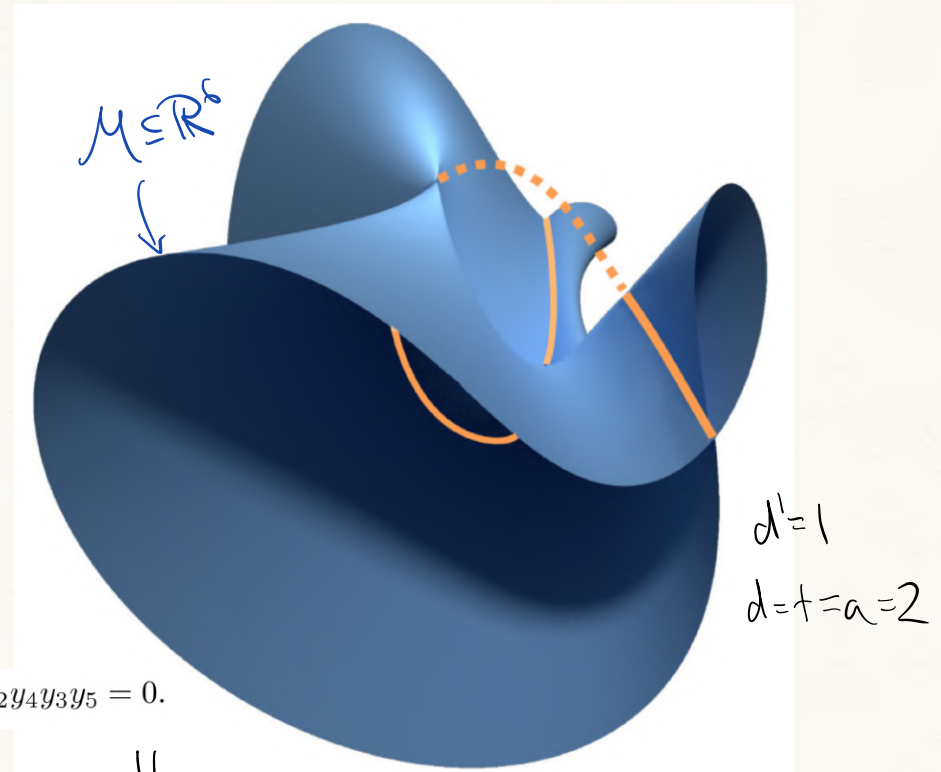| $k$ | complex critical point set |
|---|---|
| 0 | 14 points |
| 1 | 14 points |
| 2 | 4 points + a curve |
| 3 | a surface |
| 4 | a 3-dimensional subvariety |
| 5 | a 4-dimensional subvariety |

# Lightning Self-Attention (single head, single layer)

$$\mathbb{R}^{d \times t} \longrightarrow \mathbb{R}^{d' \times t}$$

$$X \longmapsto V X X^T K^T Q X$$

learnable parameters
$V \in \mathbb{R}^{d' \times d}$, $K, Q \in \mathbb{R}^{a \times d}$

$M \subseteq \mathbb{R}^6$

$d' = 1$
$d = t = a = 2$

$$y_1^2 y_6^2 + y_4^2 y_3^2 + y_1 y_3 y_5^2 + y_2^2 y_4 y_6 - 2y_1 y_4 y_3 y_6 - y_2 y_1 y_6 y_5 - y_2 y_4 y_3 y_5 = 0.$$

For almost all PD matrices $Q$,
$$\mathrm{EDD}_Q(M) = 14.$$

What happens if $Q$ becomes degenerate?

(i.e., $Q$ is symmetric positive semidefinite)

$K := \dim \ker Q$

| $k$ | complex critical point set |
|---|---|
| 0 | 14 points |
| 1 | 14 points |
| 2 | 4 points + a curve |
| 3 | a surface |
| 4 | a 3-dimensional subvariety |
| 5 | a 4-dimensional subvariety |

$M \cap (\ker(Q) + u)$
↳ zero loss solutions!

# In general

$M \subseteq \mathbb{R}^n$ algebraic variety, $d := \dim M$.

$Q$ symmetric positive semi-definite $n \times n$ matrix

$K := \ker Q$

$\pi : \mathbb{R}^n \to K^\perp$

turns $Q$ into nondegenerate quadric

# In general

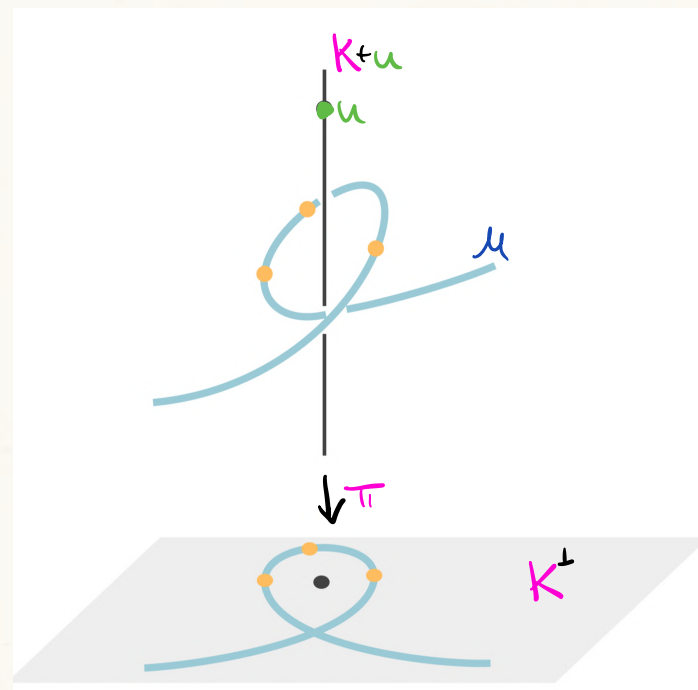$M \subseteq \mathbb{R}^n$ algebraic variety, $d := \dim M$.

$Q$ symmetric positive semi-definite $n \times n$ matrix $\quad \Big| \quad \pi : \mathbb{R}^n \to K^\perp$

$K := \ker Q$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ turns $Q$ into nondegenerate quadric

**Case 1:** Let $k < n - d$.

For almost all $Q$ with $k = \dim K$ and almost all $u \in \mathbb{R}^n$,

$EDD_Q(M)$
$\quad \Big\|$
$EDD_{\pi(Q)}(\pi(M))$
$\left\{ \begin{array}{l} \text{critical points of } \displaystyle\min_{x \in M \setminus Sing(M)} \|x - u\|_Q^2 \\[2em] \updownarrow \; 1:1 \\[2em] \text{critical points of } \displaystyle\min_{x \in \pi(M) \setminus Sing(\pi(M))} \|x - \pi(u)\|_{\pi(Q)}^2 \end{array} \right.$

# In general

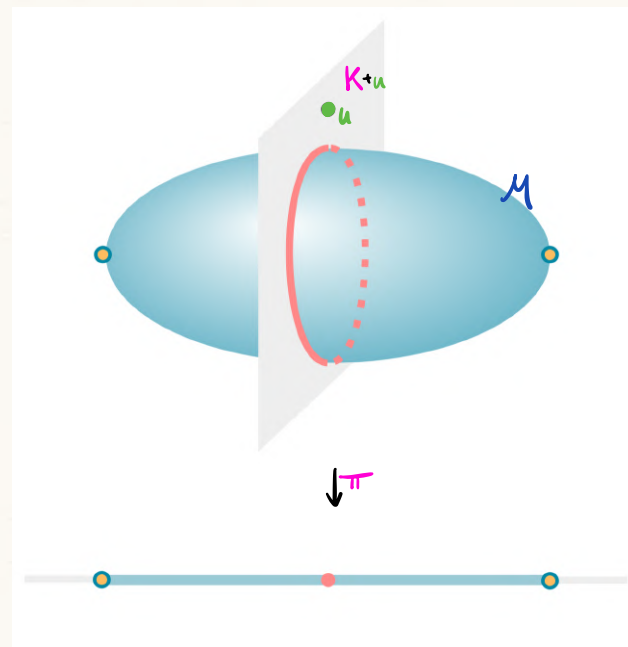$M \subseteq \mathbb{R}^n$ algebraic variety, $d := \dim M$.

$Q$ symmetric positive semi-definite $n \times n$ matrix
$K := \ker Q$

$\pi : \mathbb{R}^n \to K^\perp$
turns $Q$ into nondegenerate quadric

**Case 2:** Let $k \geq n - d$.

For almost all $Q$ with $k = \dim K$ and almost all $u \in \mathbb{R}^n$, we have
2 types of critical points of $\min\limits_{x \in M \setminus \text{Sing}(M)} \|x - u\|_Q^2$ :

Ⓐ $(K+u) \cap M$ : zero loss solutions

# In general

$M \subseteq \mathbb{R}^n$ algebraic variety, $d := \dim M$.

$Q$ symmetric positive semi-definite $n \times n$ matrix
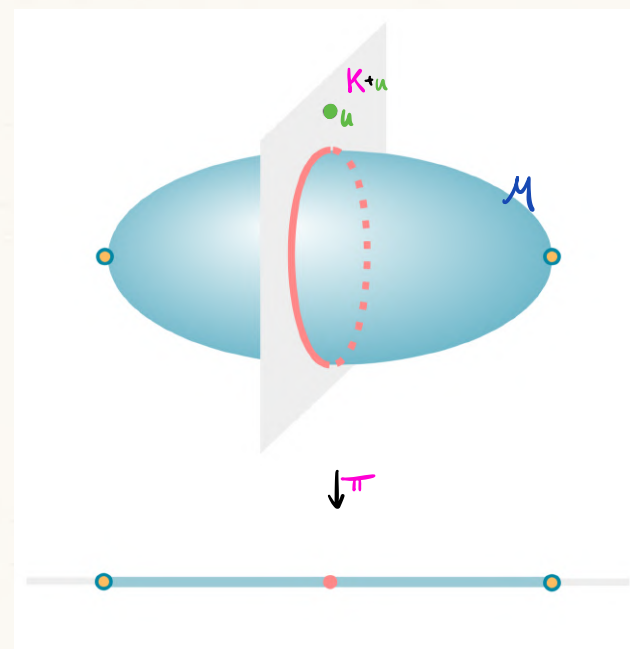$K := \ker Q$

$\pi : \mathbb{R}^n \to K^\perp$
turns $Q$ into nondegenerate quadric

**Case 2:** Let $k \geq n - d$.

For almost all $Q$ with $k = \dim K$ and almost all $u \in \mathbb{R}^n$, we have
2 types of critical points of $\min\limits_{x \in M \setminus Sing(M)} \|x - u\|^2_Q$ :

(A) $(K + u) \cap M$ : zero loss solutions

(B) finitely many on the **ramification locus** $\mathrm{Ram}(\pi|_X)$
$:= \{$critical points of $\pi|_X\}$

# In general

$M \subseteq \mathbb{R}^n$ algebraic variety, $d := \dim M$.

$Q$ symmetric positive semi-definite $n \times n$ matrix | $\pi : \mathbb{R}^n \to K^\perp$
$K := \ker Q$ | turns $Q$ into nondegenerate quadric

**Case 2:** Let $k \geq n - d$.

For almost all $Q$ with $k = \dim K$ and almost all $u \in \mathbb{R}^n$, we have 2 types of critical points of $\min\limits_{x \in M \setminus Sing(M)} \|x - u\|^2_Q$ :

(A) $(K + u) \cap M$ : zero loss solutions

(B) finitely many on the **ramification locus** $Ram(\pi|_X)$
$:= \{$critical points of $\pi|_X\}$

$\Big\uparrow 1:1$

$EDD_{\pi(Q)}(Br) \leftarrow$ critical points of $\min\limits_{x \in Br(\pi|_X)} \|x - \pi(u)\|^2_{\pi(Q)}$
$\leftarrow$ **Branch locus** $\pi(Ram)$

# In general

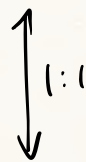$M \subseteq \mathbb{R}^n$ algebraic variety, $d := \dim M$.

$Q$ symmetric positive semi-definite $n \times n$ matrix
$K := \ker Q$

$\pi : \mathbb{R}^n \to K^\perp$
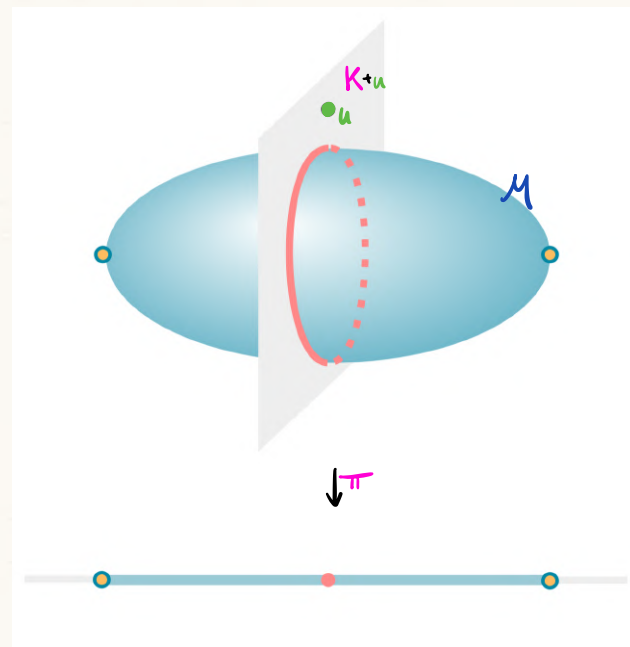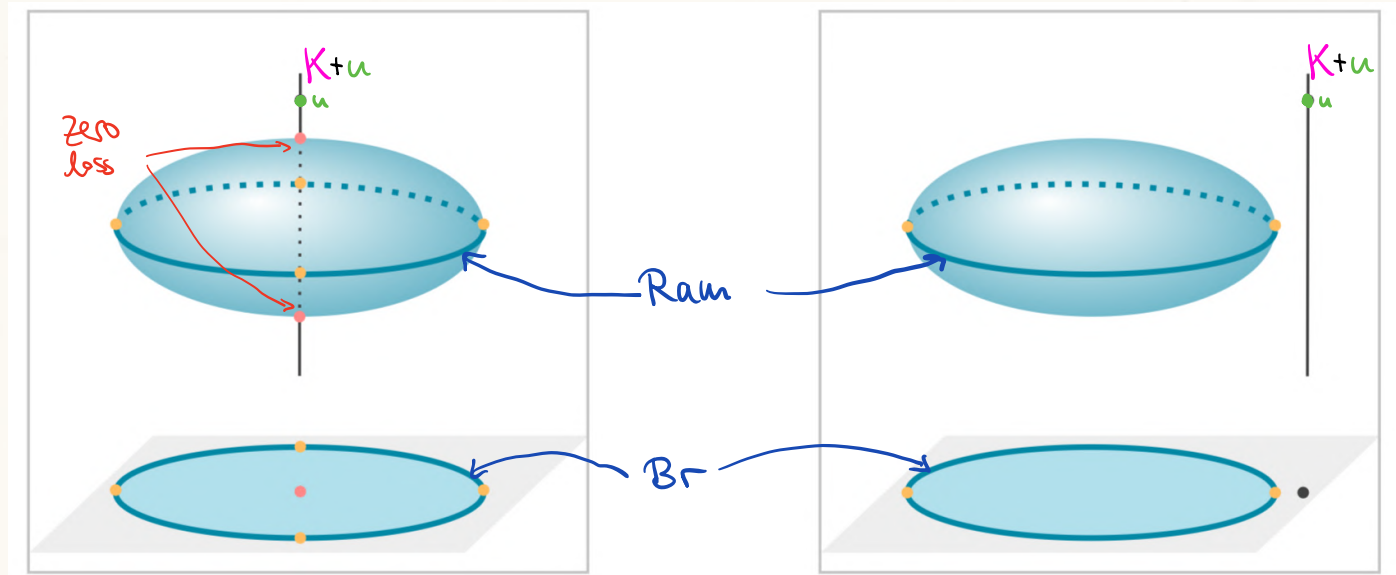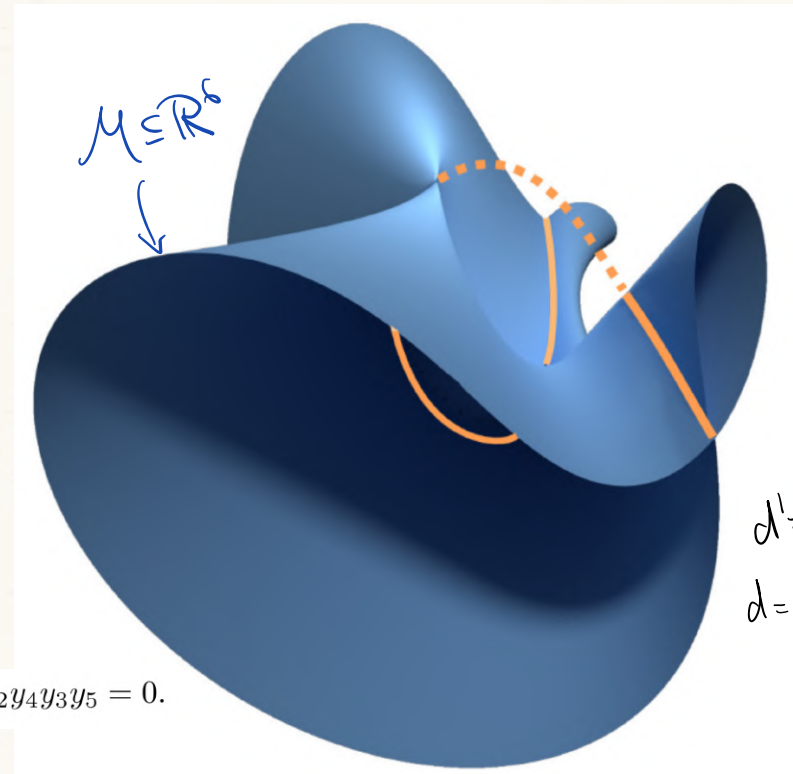turns $Q$ into nondegenerate quadric

## Case 2: let $k \geq n - d$.



Induced bias towards Ram!

depends only on $K$ (not on $Q$) & not on $u$

# Lightning Self-Attention (single head, single layer)

$$\mathbb{R}^{d \times t} \longrightarrow \mathbb{R}^{d' \times t}$$

$$X \longmapsto V X X^T K^T Q X$$

learnable parameters
$V \in \mathbb{R}^{d' \times d}$, $K, Q \in \mathbb{R}^{a \times d}$

$M \subseteq \mathbb{R}^6$

$d' = 1$
$d = t = a = 2$



$$y_1^2 y_6^2 + y_4^2 y_3^2 + y_1 y_3 y_5^2 + y_2^2 y_4 y_6 - 2 y_1 y_4 y_3 y_6 - y_2 y_1 y_6 y_5 - y_2 y_4 y_3 y_5 = 0.$$

---

$Q =$ from MSE loss with general dataset $S$

| $|S|$ | $k = \dim K$ | complex critical point set |
|---|---|---|
| $\geq 3$ | 0 | 14 points |
| 2 | 2 | a curve and two lines |
| 1 | 4 | a 3-dimensional subvariety |

$Q =$ general symmetric positive semidefinite

| $k$ | complex critical point set |
|---|---|
| 0 | 14 points |
| 1 | 14 points |
| 2 | 4 points + a curve |
| 3 | a surface |
| 4 | a 3-dimensional subvariety |
| 5 | a 4-dimensional subvariety |

$K := \dim \ker Q$

$M \cap (\ker(Q) + u)$
$\hookrightarrow$ zero loss solutions!

# algebraic neural network theory – an emerging field

Kileel, Trager, Bruna: On the expressive power of deep polynomial neural networks. **NeurIPS** 2019.

Trager, Kohn, Bruna: Pure and spurious critical points: a geometric study of linear networks. **ICLR** 2020.

Kohn, Merkh, Montúfar, Trager: Geometry of linear convolutional networks. **SIAM Journal on Applied Algebra & Geometry** 2022.

Kohn, Montúfar, Shahverdi, Trager: Function space & critical points of linear CNNs. **SIAM Journal on Applied Algebra & Geometry** 2024.

Kubjas, Li, Wiesmann: Geometry of polynomial neural networks. **Algebraic Statistics** 2024.

Marchetti, Shahverdi, Mereta, Trager, Kohn: Position: Algebra unveils deep learning – An invitation to Neuroalgebraic Geometry. **ICML** 2025.

Mody, Zubkov: Geometry of Rank Constraints in Shallow Polynomial Neural Networks. **ICML** 2025 Workshop.

Zhang, Kileel: Covering number of real algebraic varieties and beyond: Improved bounds and applications. **FoCM** 2025.

Finkel, Rodriguez, Wu, Yahl: Activation thresholds and expressiveness of polynomial neural networks. **Algebraic Statistics** 2025.

Arjevani, Bruna, Kileel, Polak, Trager: Geometry and optimization of shallow polynomial networks. arXiv:2501.06074.

Shahverdi, Marchetti, Kohn: On the geometry and optimization of polynomial convolutional networks. **AISTATS** 2025.

Henry, Marchetti, Kohn: Geometry of lightning self-attention: Identifiability and dimension. **ICLR** 2025.

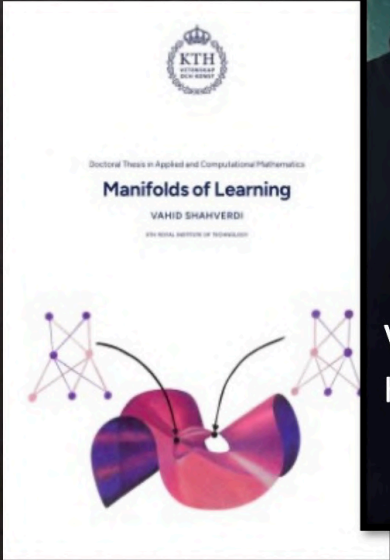Massarenti, Mella: The Alexander-Hirschowitz theorem for neurovarieties. arXiv:2511.19703.

Grosdos, Robeva, Zubkov: Algebraic geometry of rational neural networks. arXiv:2509.11088.

Shahverdi, Marchetti, Kohn: Learning on a razor's edge: the singularity bias of polynomial neural networks. ~~arXiv:2505.11846.~~ *ICLR 2026*

Shahverdi: Algebraic complexity and neurovariety of linear convolutional networks. **Acta Univ. Sapientiae Math.** 2025.
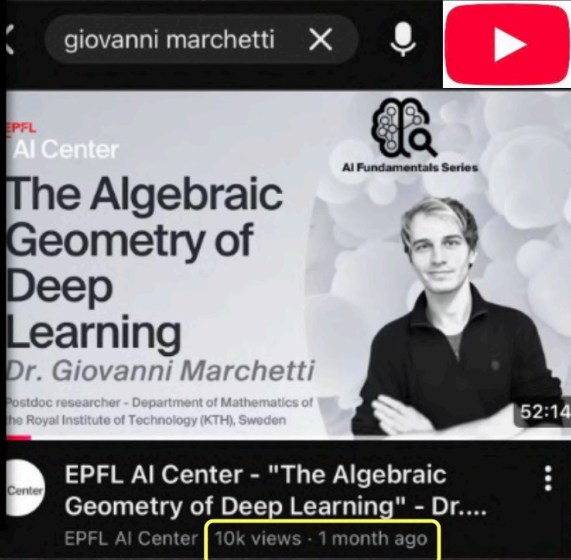
Usevich, Dérand, Borsoi, Clausel: Identifiability of Deep Polynomial Neural Networks. **NeurIPS Oral** 2025.

special thanks to



Vahid Shahverdi
KTH

Giovanni Marchetti
KTH

Matthew Trager
AWS AI Labs, NY

Stefano Mereta
CUNEF Madrid

Nathan Henry
UC Berkeley

Guido Montúfar
UCLA & MPI MiS Leipzig

Joan Bruna
Courant Institute, NYU

Paul Breiding & Erin Connelly
Univ. Osnabrück