

Algebra & Geometry of Neural Networks

NTK approach

neural
target
kernel

target
network

increase
width

∞ -width
network

linearized models
of ∞ dimension

neural
target
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

Stone - Weierstraß

continuous
functions
↙

Let X compact Hausdorff space & A subalgebra of $C(X, \mathbb{R})$ containing a nonzero constant function.

A is dense in $C(X, \mathbb{R})$
in supremum norm

\Leftrightarrow A separates points
(i.e., $\forall x \neq y \in X \exists f \in A: f(x) \neq f(y)$)

Cor.: $X \subseteq \mathbb{R}^n$ compact, $f: X \rightarrow \mathbb{R}^m$ continuous, $\varepsilon > 0$.

$\Rightarrow \exists p: X \rightarrow \mathbb{R}^m$ polynomial function such that
 $\forall x \in X: \|f(x) - p(x)\| < \varepsilon$.

Example: MLPs multilayer perceptrons

$$\alpha_L \circ \sigma \circ \dots \circ \sigma \circ \alpha_2 \circ \sigma \circ \alpha_1$$

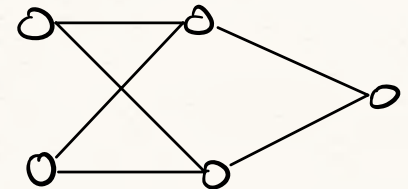
$\alpha_i =$ learnable affine linear functions

$\sigma =$ nonlinear activation function, applied entrywise

We assume: σ is a univariate **polynomial**

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

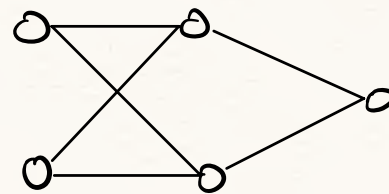
$$[e \ f] \sigma \left(\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \right)$$



Which functions does this MLP parametrize?

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?

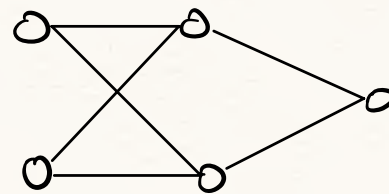
$$\begin{aligned} & e(ax+by)^2 + f(cx+dy)^2 \\ &= \underbrace{(a^2e + c^2f)}_A x^2 + \underbrace{2(ab e + cd f)}_B xy + \underbrace{(b^2e + d^2f)}_C y^2 \end{aligned}$$

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?

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?

$$\begin{aligned} & e(ax+by)^2 + f(cx+dy)^2 \\ &= \underbrace{(a^2e + c^2f)}_A x^2 + \underbrace{2(ab e + cd f)}_B xy + \underbrace{(b^2e + d^2f)}_C y^2 \end{aligned}$$

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

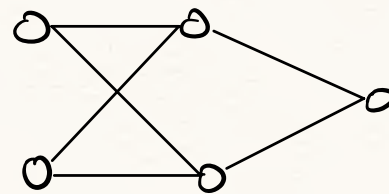
← homogeneous quadratic polynomials in x,y
i.e., are all values for A, B, C possible?

YES

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

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

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



Which functions does this MLP parametrize?

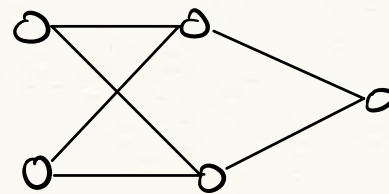
$$\begin{aligned} & e(ax+by)^3 + f(cx+dy)^3 \\ &= \underbrace{(a^3e + c^3f)}_A x^3 + \underbrace{3(a^2be + c^2df)}_B x^2y + \underbrace{3(ab^2e + cd^2f)}_C xy^2 + \underbrace{(b^3e + d^3f)}_D y^3 \end{aligned}$$

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?

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?

$$\begin{aligned} & e(ax+by)^3 + f(cx+dy)^3 \\ &= \underbrace{(a^3e + c^3f)}_A x^3 + \underbrace{3(a^2be + c^2df)}_B x^2y + \underbrace{3(ab^2e + cd^2f)}_C xy^2 + \underbrace{(b^3e + d^3f)}_D y^3 \end{aligned}$$

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.

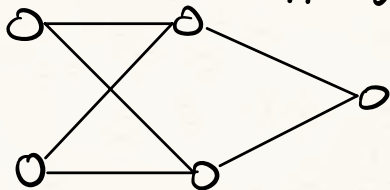
| | | |
|-----|-----|----|
| A | $=$ | 1 |
| B | $=$ | 0 |
| C | $=$ | -1 |
| D | $=$ | 0 |

Macaulay 2

Neuromanifolds

A **parametric machine learning** model is a map $\mu: \Theta \times X \rightarrow Y$.
parameters \uparrow \uparrow \uparrow
inputs outputs

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

Examples:  no bias

$$\sigma(x) = x^2$$

$$\Rightarrow \mathcal{M} = \mathbb{R}[x, y]_2$$

$$\sigma(x) = x^3$$

$$\rightarrow \mathcal{M} \subsetneq \mathbb{R}[x, y]_3$$

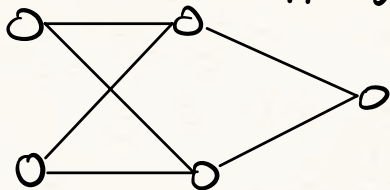
$$\sigma(x) = x$$

$$\Rightarrow ?$$

Neuromanifolds

A **parametric machine learning** model is a map $\mu: \Theta \times X \rightarrow Y$.
 parameters \uparrow \uparrow \uparrow
 Θ X Y
 inputs outputs

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

Examples:  no bias

$$\sigma(x) = x^2$$

$$\Rightarrow \mathcal{M} = \mathbb{R}[x, y]_2$$

$$\sigma(x) = x^3$$

$$\Rightarrow \mathcal{M} \subsetneq \mathbb{R}[x, y]_3$$

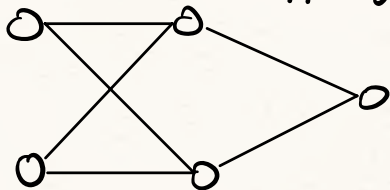
$$\sigma(x) = x$$

$$\Rightarrow \mathcal{M} = \mathbb{R}^{1 \times 2}$$

Neuromanifolds

A **parametric machine learning** model is a map $\mu: \Theta \times X \rightarrow Y$.
 parameters \uparrow \uparrow \uparrow
 Θ X Y
 inputs outputs

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

Examples:  no bias

$$\sigma(x) = x^2$$

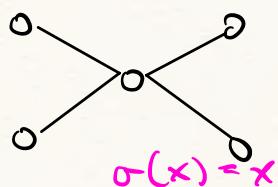
$$\Rightarrow \mathcal{M} = \mathbb{R}[x, y]_2$$

$$\sigma(x) = x^3$$

$$\Rightarrow \mathcal{M} \subsetneq \mathbb{R}[x, y]_3$$

$$\sigma(x) = x$$

$$\Rightarrow \mathcal{M} = \mathbb{R}^{1 \times 2}$$



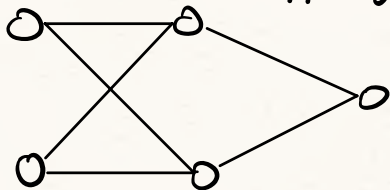
$$\begin{bmatrix} c \\ a \end{bmatrix} \begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\Rightarrow \mathcal{M} = ?$$

Neuromanifolds

A **parametric machine learning** model is a map $\mu: \Theta \times X \rightarrow Y$.
 parameters \uparrow \uparrow \uparrow
 Θ X Y
 inputs outputs

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

Examples:  no bias

$$\sigma(x) = x^2$$

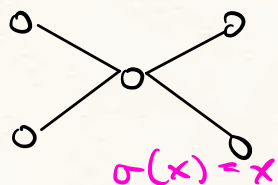
$$\Rightarrow \mathcal{M} = \mathbb{R}[x, y]_2$$

$$\sigma(x) = x^3$$

$$\Rightarrow \mathcal{M} \subsetneq \mathbb{R}[x, y]_3$$

$$\sigma(x) = x$$

$$\Rightarrow \mathcal{M} = \mathbb{R}^{1 \times 2}$$



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

$$\Rightarrow \mathcal{M} = \{ W \in \mathbb{R}^{2 \times 2} \mid \text{rk}(W) \leq 1 \}$$

Linear MLPs:

$\alpha_L \circ \dots \circ \alpha_2 \circ \alpha_1$, where
 $\alpha_i: \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$ linear

$\Rightarrow M = 2$

Linear MLPs: $\alpha_L \circ \dots \circ \alpha_2 \circ \alpha_1$, where
 $\alpha_i: \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$ linear

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

Linear MLPs: $\alpha_L \circ \dots \circ \alpha_2 \circ \alpha_1$, where
 $\alpha_i: \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$ 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: $\alpha_L \circ \sigma \circ \dots \circ \sigma \circ \alpha_2 \circ \sigma \circ \alpha_1$, where
 $\alpha_i: \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$ affine linear

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

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

Linear MLPs: $\alpha_L \circ \dots \circ \alpha_2 \circ \alpha_1$, where
 $\alpha_i: \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$ 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: $\alpha_L \circ \sigma \circ \dots \circ \sigma \circ \alpha_2 \circ \sigma \circ \alpha_1$, where
 $\alpha_i: \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$ affine linear

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

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

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

Linear MLPs: $\alpha_L \circ \dots \circ \alpha_2 \circ \alpha_1$, where
 $\alpha_i: \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$ 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: $\alpha_L \circ \sigma \circ \dots \circ \sigma \circ \alpha_2 \circ \sigma \circ \alpha_1$, where
 $\alpha_i: \mathbb{R}^{d_{i-1}} \rightarrow \mathbb{R}^{d_i}$ affine linear

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

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

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

Polynomial MLPs are the only ones with that property!

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

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

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

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

Theorem 1:

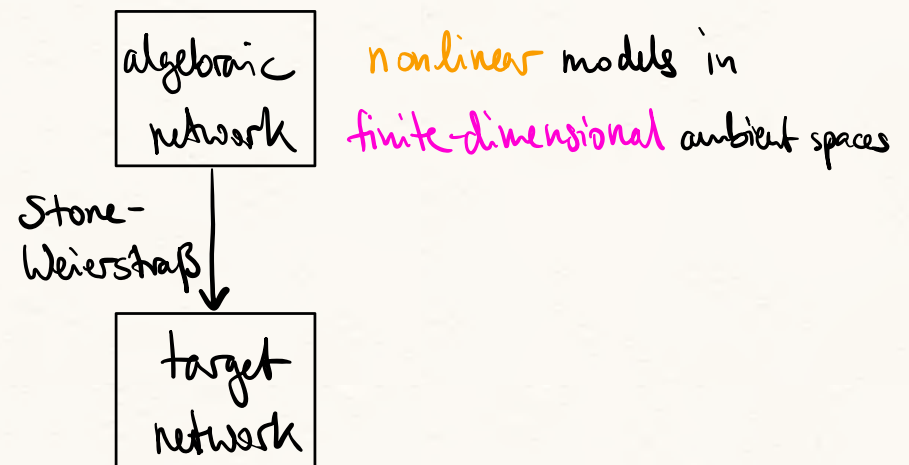
Let $\sigma \in M$. Set

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

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

polynomials are the choice
to approximate networks with
finite-dimensional models

AG approach



Network training = 'distance' minimization

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

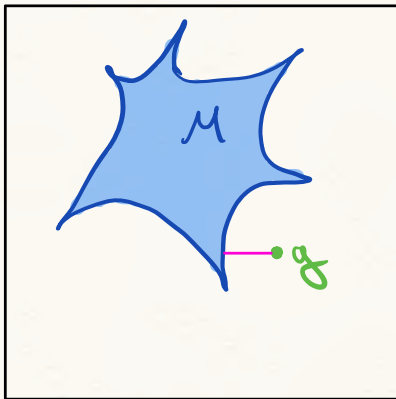
$S \subseteq \mathbb{R}^{d_0} \times \mathbb{R}^{d_L}$ finite dataset,

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

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

Proposition: There is a pseudometric $\text{dist}: V \times V \rightarrow \mathbb{R}_{\geq 0}$ 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



Why?

Network training = 'distance' minimization

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

\nwarrow neuromanifold

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

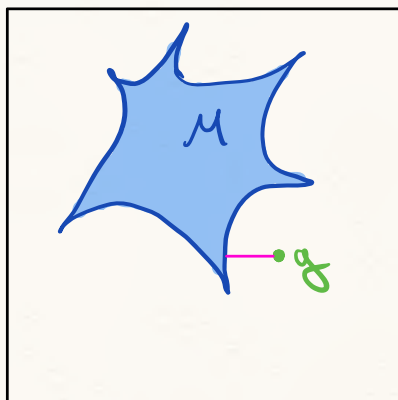
\nwarrow mean squared error

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

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

Proposition: There is a pseudometric $\text{dist}: V \times V \rightarrow \mathbb{R}_{\geq 0}$ 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



Assume: $d_L = 1$

Let $v_D: (x_1, \dots, x_{d_0}) \mapsto (\text{all monomials in } x_1, \dots, x_{d_0} \text{ of degree } \leq \mathbb{D})$,
 c_f be coefficient vector of $f \in V$ such that $f(x) = v_D(x) \cdot c_f$,

Veronese embedding ↘

Network training = 'distance' minimization

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

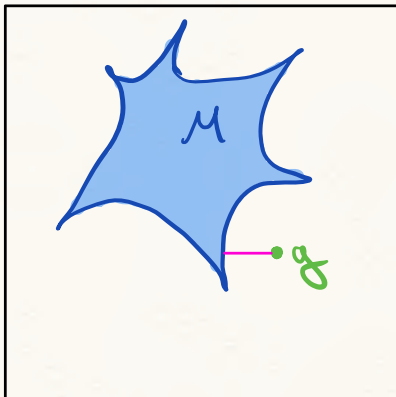
$S \subseteq \mathbb{R}^{d_0} \times \mathbb{R}^{d_L}$ finite dataset,

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

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

Proposition: There is a pseudometric $\text{dist}: V \times V \rightarrow \mathbb{R}_{\geq 0}$ 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 $f \in \mathcal{M}$.

\vee



Assume: $d_L = 1$

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

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

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

Veronese embedding \curvearrowright

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

Network training = 'distance' minimization

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

\nwarrow *neuromanifold*

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

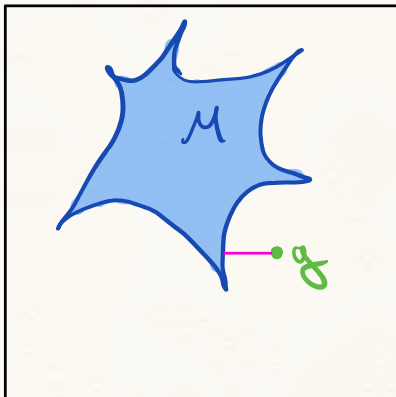
$$\text{MSE loss: } \mathcal{L}(f) := \sum_{(a,b) \in S} \|f(a) - b\|^2$$

\nwarrow *mean squared error*

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

Proposition: There is a pseudometric $\text{dist}: V \times V \rightarrow \mathbb{R}_{\geq 0}$ 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



Assume: $d_L = 1$

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

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

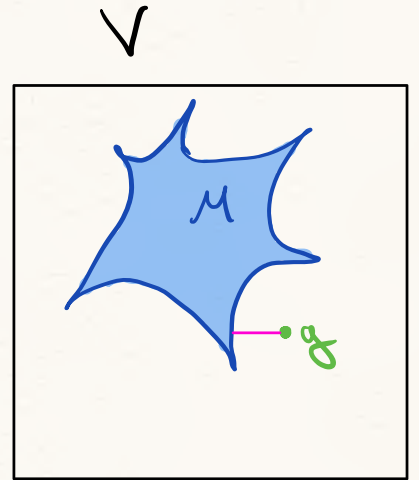
A & B matrices whose rows are $v_D(a)$ & b , resp., over all $(a,b) \in S$

Veronese embedding \curvearrowright

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

\nwarrow *pseudoinverse*
 \nwarrow $\|c\|_Q := c^T Q c$

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



Observations ($d_L=1$):

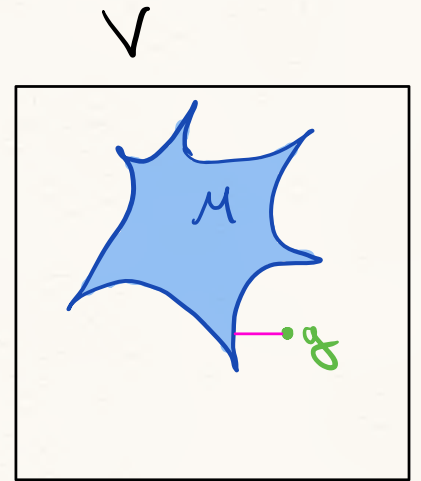
① $A^T A$ depends only on input data,
 $A^* B$ on both input & output

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

③

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

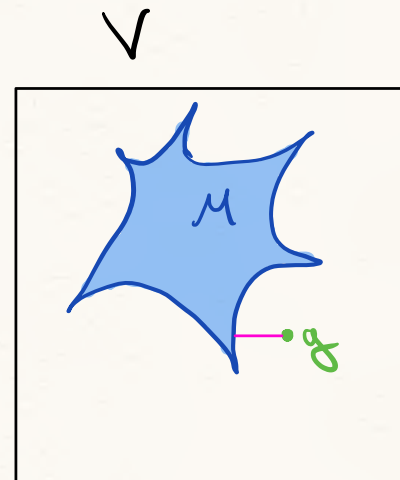
$$\arg\min_{f \in \mathcal{M}} L(f) = \arg\min_{f \in \mathcal{M}} \|C_f - A^+ B\|_{A^T A}^2$$



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 $|S| < \dim V \rightarrow$ pseudometric
(LLMs: $|S| < \dim M$)
- ③ even when $|S| \gg \dim V$, $A^T A$ is not an arbitrary symmetric PD matrix,
 while $A^+ B$ yields all vectors $\in \mathbb{R}^{\dim V}$
Why?
Which matrices can be obtained?
 (try for $d_0=1$: $v(x) = (1, x, x^2, \dots, x^D)$)

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

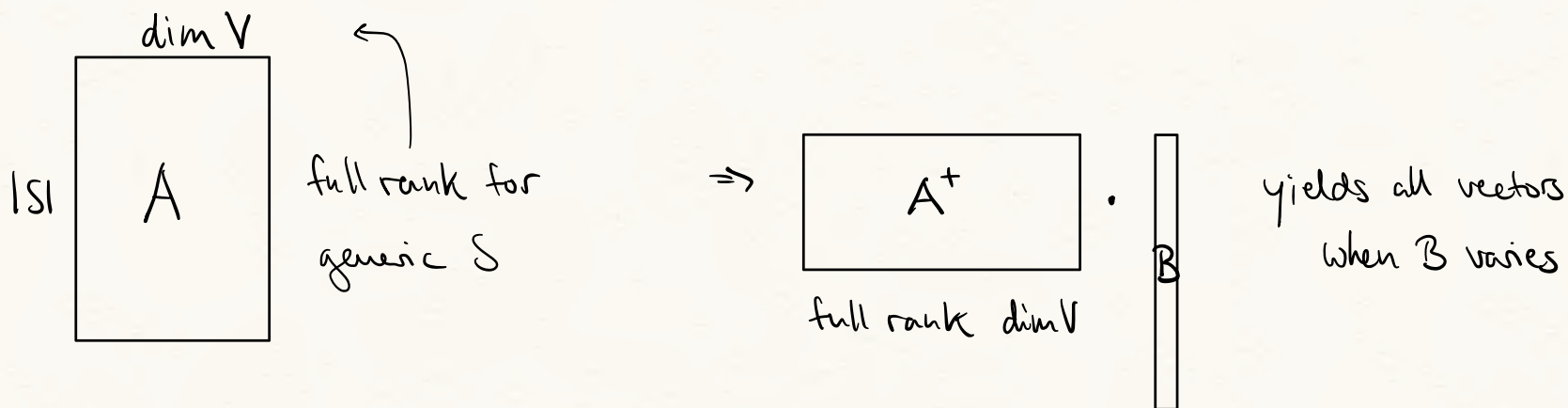


Observations ($d_L=1$):

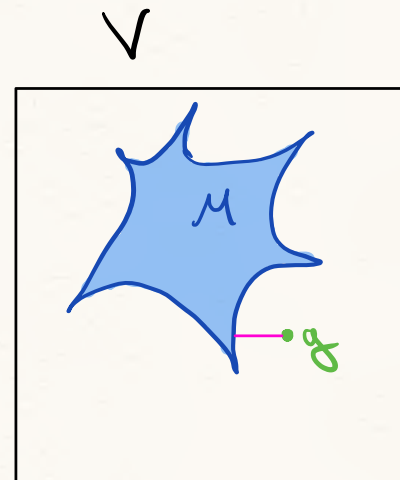
① $A^T A$ depends only on input data,
 $A^+ B$ on both input & output

② $A^T A \in \mathbb{R}^{\dim V \times \dim V}$ is rank-deficient whenever $|S| < \dim V \rightarrow$ pseudometric (LLMs: $|S| < \dim M$)

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



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



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 $|S| < \dim V \rightarrow$ pseudometric (LLMs: $|S| < \dim M$)

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

$$A^T A = \begin{matrix} & i \rightarrow \\ \begin{array}{|c|c|c|} \hline | & & | \\ \hline v(a_1) & \dots & v(a_{|S|}) \\ \hline | & & | \\ \hline \end{array} & \begin{array}{|c|} \hline v(a_1) \\ \hline \vdots \\ \hline v(a_{|S|}) \\ \hline \end{array} \end{matrix}$$

has (i,j) entry $\sum_{(a,b) \in S} \underbrace{v_i(a) v_j(a)}_{\text{monomial of degree } \leq 2D}$
 that can be factored in several ways

Ex.: $d_0 = 1$

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

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

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

Ex.: $d_0 = 1$

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

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

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

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

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

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

Network training = 'distance' minimization

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

\nwarrow neuromanifold

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

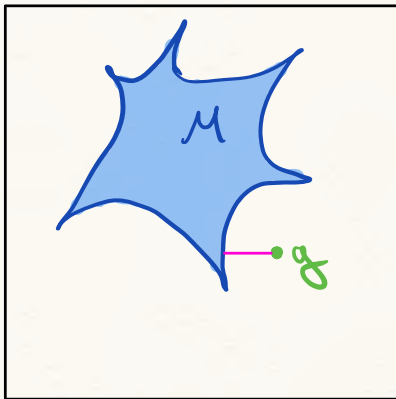
\nwarrow mean squared error

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

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

Proposition: There is a pseudometric $\text{dist}: V \times V \rightarrow \mathbb{R}_{\geq 0}$ 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



$d_L > 1$

$$f = (f_1, \dots, f_{d_L}), \quad C_f := \begin{bmatrix} | & & | \\ c_{f_1} & \dots & c_{f_{d_L}} \\ | & & | \end{bmatrix}$$

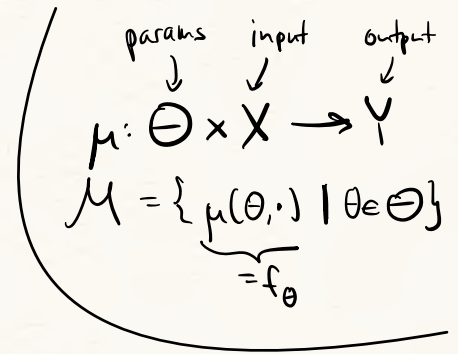
$$\Rightarrow f(x) = v_D(x) \cdot C_f$$

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

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

Loss Landscape

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



Loss Landscape

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

$$\begin{array}{ccc} \text{params} & \text{input} & \text{output} \\ \downarrow & \downarrow & \downarrow \\ \mu: \Theta \times X & \rightarrow & Y \\ \mathcal{M} = \{ \underbrace{\mu(\theta, \cdot)}_{=f_\theta} \mid \theta \in \Theta \} \end{array}$$

can be studied in a decoupled way:

$$\begin{array}{ccccc} \Theta & \xrightarrow{\quad} & \mathcal{M} & \xrightarrow{\mathcal{L}} & \mathbb{R} \\ \theta & \mapsto & f_\theta & & \end{array}$$



loss landscape in function space:

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

Loss Landscape

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

$$\begin{array}{ccc} \text{params} & \text{input} & \text{output} \\ \downarrow & \downarrow & \downarrow \\ \mu: \Theta \times X & \rightarrow & Y \\ \mathcal{M} = \{ \underbrace{\mu(\theta, \cdot)}_{=f_\theta} \mid \theta \in \Theta \} \end{array}$$

can be studied in a decoupled way:

$$\begin{array}{ccccc} \Theta & \xrightarrow{\quad} & \mathcal{M} & \xrightarrow{\mathcal{L}} & \mathbb{R} \\ \theta & \mapsto & f_\theta & & \end{array}$$



loss landscape in function space:

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

How?

Geometry of \mathcal{M} affects loss landscape!

Which geometric properties does \mathcal{M} have?