# understanding Linear Convolutional Neural Networks via sparse factorizations of real polynomials

Kathlén Kohn



joint work with

Guido Montúfar          Vahid Shahverdi          Matthew Trager
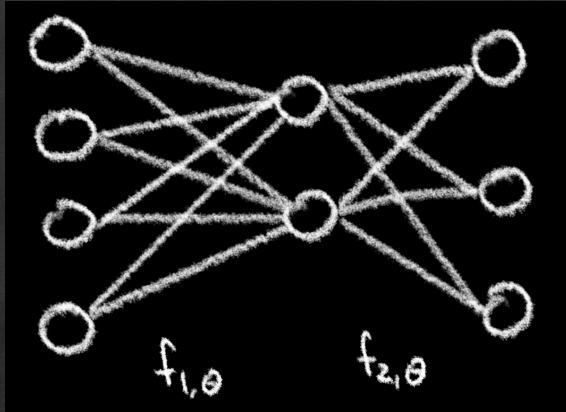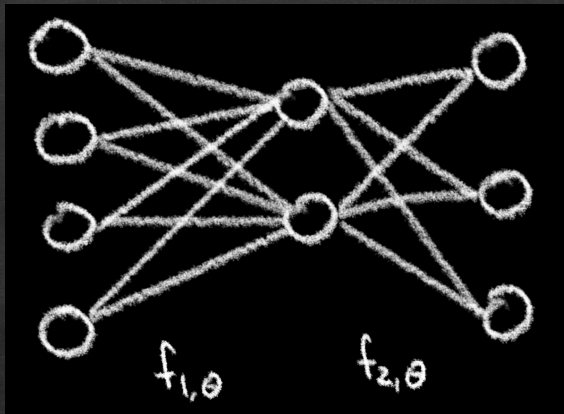MPI MiS Leipzig & UCLA          KTH          Amazon Alexa AI, NYC

# feedforward neural networks

# feedforward neural networks



are parametrized families of functions

$$\mu : \mathbb{R}^N \longrightarrow \mathcal{M},$$
$$\theta \longmapsto f_{L,\theta} \circ \ldots \circ f_{1,\theta}$$
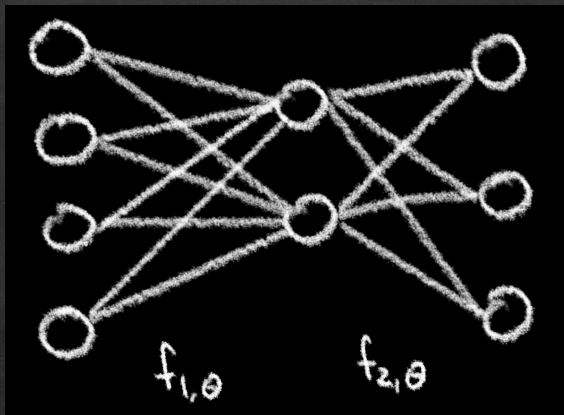
# feedforward neural networks



are parametrized families of functions
$$\mu : \mathbb{R}^N \longrightarrow \mathcal{M},$$
$$\theta \longmapsto f_{L,\theta} \circ \ldots \circ f_{1,\theta}$$
$\mathcal{M} =$ function space / neuromanifold, $L = \#$ layers

# training a network

Given training data $\mathcal{D}$, the goal is to minimize the loss

$$\mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}.$$

# training a network

Given training data $\mathcal{D}$, the goal is to minimize the loss

$$\mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}.$$

**Geometric questions:**

◆ How does the network architecture affect the geometry of the function space?

◆ How does the geometry of the function space impact the training of the network?

# training a network

Given training data $\mathcal{D}$, the goal is to minimize the loss

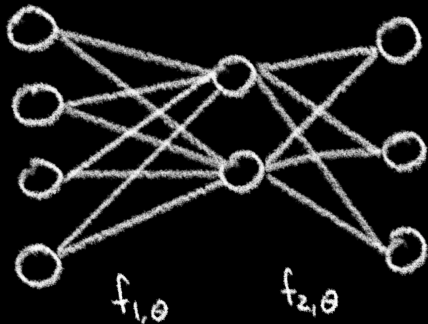$$\mathbb{R}^N \xrightarrow{\ \mu\ } \mathcal{M} \xrightarrow{\ \ell_{\mathcal{D}}\ } \mathbb{R}.$$

**Geometric questions:**
- How does the network architecture affect the geometry of the function space?
- How does the geometry of the function space impact the training of the network?

**In this talk:**
What is the impact of changing from dense layers to convolutional layers?

# linear dense networks



In this example:

$$\mu : \mathbb{R}^{2\times 4} \times \mathbb{R}^{3\times 2} \longrightarrow \mathbb{R}^{3\times 4},$$
$$(W_1, W_2) \longmapsto W_2 W_1.$$

# linear dense networks



In this example:

$$\mu : \mathbb{R}^{2\times 4} \times \mathbb{R}^{3\times 2} \longrightarrow \mathbb{R}^{3\times 4},$$
$$(W_1, W_2) \longmapsto W_2 W_1.$$

$$\mathcal{M} = \{W \in \mathbb{R}^{3\times 4} \mid \operatorname{rank}(W) \leq 2\}$$
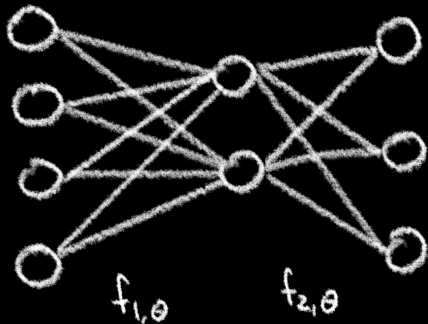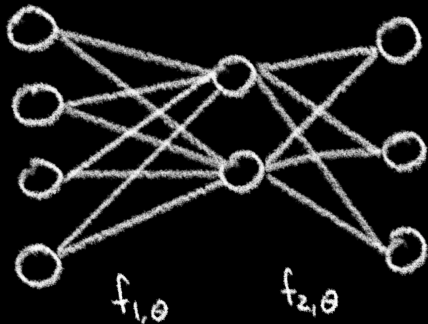
# linear dense networks



In this example:

$$\mu : \mathbb{R}^{2 \times 4} \times \mathbb{R}^{3 \times 2} \longrightarrow \mathbb{R}^{3 \times 4},$$
$$(W_1, W_2) \longmapsto W_2 W_1.$$

$$\mathcal{M} = \{W \in \mathbb{R}^{3 \times 4} \mid \mathrm{rank}(W) \leq 2\}$$

In general:

$$\mu : \mathbb{R}^{k_1 \times k_0} \times \mathbb{R}^{k_2 \times k_1} \times \ldots \times \mathbb{R}^{k_L \times k_{L-1}} \longrightarrow \mathbb{R}^{k_L \times k_0},$$
$$(W_1, W_2, \ldots, W_L) \longmapsto W_L \cdots W_2 W_1.$$

$\mathcal{M} = \{W \in \mathbb{R}^{k_L \times k_0} \mid \mathrm{rank}(W) \leq \min(k_0, \ldots, k_L)\}$ is an algebraic variety and we know its singularities etc.

# Linear Convolutional Networks (LCNs)
### with 1D-convolutions

# Linear Convolutional Networks (LCNs)

## with 1D-convolutions



$$\mu : \mathbb{R}^3 \times \mathbb{R}^2 \longrightarrow \mathbb{R}^5,$$

$$(u, v) \longmapsto T_{v,1} \, T_{u,2}, \text{ where}$$

$$T_{u,2} = \left[ \begin{array}{ccccc} u_0 & u_1 & u_2 & 0 & 0 \\ 0 & 0 & u_0 & u_1 & u_2 \end{array} \right]$$

$$T_{v,1} = \left[ \begin{array}{cc} v_0 & v_1 \end{array} \right]$$

# Linear Convolutional Networks (LCNs)
## with 1D-convolutions



$$\mu : \mathbb{R}^3 \times \mathbb{R}^2 \longrightarrow \mathbb{R}^5,$$
$$(u, v) \longmapsto T_{v,1}\, T_{u,2}, \text{ where}$$

$$T_{u,2} = \begin{bmatrix} u_0 & u_1 & u_2 & 0 & 0 \\ 0 & 0 & u_0 & u_1 & u_2 \end{bmatrix}$$

$$T_{v,1} = \begin{bmatrix} v_0 & v_1 \end{bmatrix}$$

In general: $\mu : (w_1, \ldots, w_L) \mapsto T_{w_L, s_L} \cdots T_{w_1, s_1}$, where

$$T_{w,s} = \begin{bmatrix} w_0 & \cdots & w_s & \cdots & w_{k-1} & & & \\ & w_0 & & \cdots & & w_{k-1} & & \\ & & \ddots & & & & \ddots & \\ & & & w_0 & & \cdots & & w_{k-1} \end{bmatrix}$$

# Linear Convolutional Networks (LCNs)
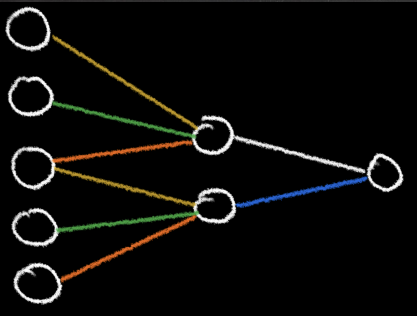## with 1D-convolutions



$$\mu : \mathbb{R}^3 \times \mathbb{R}^2 \longrightarrow \mathbb{R}^5,$$
$$(u, v) \longmapsto T_{v,1}\, T_{u,2}, \text{ where}$$

$$T_{u,2} = \begin{bmatrix} u_0 & u_1 & u_2 & 0 & 0 \\ 0 & 0 & u_0 & u_1 & u_2 \end{bmatrix}$$

$$T_{v,1} = \begin{bmatrix} v_0 & v_1 \end{bmatrix}$$

In general: $\mu : (w_1, \ldots, w_L) \mapsto T_{w_L, s_L} \cdots T_{w_1, s_1}$, where

$$T_{w,s} = \begin{bmatrix} w_0 & \cdots & w_s & \cdots & w_{k-1} & & & \\ & w_0 & & \cdots & & w_{k-1} & & \\ & & \ddots & & & & \ddots & \\ & & & w_0 & & \cdots & & w_{k-1} \end{bmatrix}$$

is a convolutional matrix of stride $s$ with filter $w$

IV - XVII

# LCNs & sparse polynomial factorization

**Observation:** $\mu(w_1, \ldots, w_L) = T_{w_L, s_L} \cdots T_{w_1, s_1}$ is again a convolutional matrix of stride $s_1 \cdots s_L$.

# LCNs & sparse polynomial factorization

**Observation:** $\mu(w_1, \ldots, w_L) = T_{w_L, s_L} \cdots T_{w_1, s_1}$ is again a convolutional matrix of stride $s_1 \cdots s_L$. Its filter can be computed via polynomial multiplication:

# LCNs & sparse polynomial factorization

**Observation:** $\mu(w_1, \ldots, w_L) = T_{w_L, s_L} \cdots T_{w_1, s_1}$ is again a convolutional matrix of stride $s_1 \cdots s_L$. Its filter can be computed via polynomial multiplication:

For $S \in \mathbb{Z}_{>0}$, let

$$\pi_S : \mathbb{R}^k \longrightarrow \mathbb{R}[x^S]_{\leq k-1},$$
$$v \longmapsto v_0 x^{S(k-1)} + v_1 x^{S(k-2)} + \ldots + v_{k-2} x^S + v_{k-1}$$

# LCNs & sparse polynomial factorization

**Observation:** $\mu(w_1, \ldots, w_L) = T_{w_L, s_L} \cdots T_{w_1, s_1}$ is again a convolutional matrix of stride $s_1 \cdots s_L$. Its filter can be computed via polynomial multiplication:

For $S \in \mathbb{Z}_{>0}$, let

$$\pi_S : \mathbb{R}^k \longrightarrow \mathbb{R}[x^S]_{\leq k-1},$$
$$v \longmapsto v_0 x^{S(k-1)} + v_1 x^{S(k-2)} + \ldots + v_{k-2} x^S + v_{k-1}$$

and $\pi_S(T_{w,s}) := \pi_S(w)$. Then:

# LCNs & sparse polynomial factorization

**Observation:** $\mu(w_1, \ldots, w_L) = T_{w_L, s_L} \cdots T_{w_1, s_1}$ is again a convolutional matrix of stride $s_1 \cdots s_L$. Its filter can be computed via polynomial multiplication:

For $S \in \mathbb{Z}_{>0}$, let

$$\pi_S : \mathbb{R}^k \longrightarrow \mathbb{R}[x^S]_{\leq k-1},$$
$$v \longmapsto v_0 x^{S(k-1)} + v_1 x^{S(k-2)} + \ldots + v_{k-2} x^S + v_{k-1}$$

and $\pi_S(T_{w,s}) := \pi_S(w)$. Then:

$$\pi_1(\mu(w_1, \ldots, w_L)) = \pi_{S_L}(w_L) \cdots \pi_{S_1}(w_1), \text{ where } S_i := s_1 \cdots s_{i-1}.$$

# LCNs & sparse polynomial factorization

**Observation:** $\mu(w_1, \ldots, w_L) = T_{w_L, s_L} \cdots T_{w_1, s_1}$ is again a convolutional matrix of stride $s_1 \cdots s_L$. Its filter can be computed via polynomial multiplication:

For $S \in \mathbb{Z}_{>0}$, let

$$\pi_S : \mathbb{R}^k \longrightarrow \mathbb{R}[x^S]_{\leq k-1},$$
$$v \longmapsto v_0 x^{S(k-1)} + v_1 x^{S(k-2)} + \ldots + v_{k-2} x^S + v_{k-1}$$

and $\pi_S(T_{w,s}) := \pi_S(w)$. Then:

$$\pi_1(\mu(w_1, \ldots, w_L)) = \pi_{S_L}(w_L) \cdots \pi_{S_1}(w_1), \text{ where } S_i := s_1 \cdots s_{i-1}.$$

Hence, we reinterpret $\mu$ as

$$\mu : \mathbb{R}[x^{S_1}]_{\leq d_1} \times \ldots \times \mathbb{R}[x^{S_L}]_{\leq d_L} \longrightarrow \mathbb{R}[x]_{\leq d_1 S_1 + \ldots + d_L S_L},$$
$$(P_1, \ldots, P_L) \longmapsto P_L \cdots P_1$$

# LCN function spaces

$$\mu : \mathbb{R}[x^{S_1}]_{\leq d_1} \times \ldots \times \mathbb{R}[x^{S_L}]_{\leq d_L} \longrightarrow \mathbb{R}[x]_{\leq d}, \text{ where } d := \sum_i d_i S_i$$
$$(P_1, \ldots, P_L) \longmapsto P_L \cdots P_1,$$

# LCN function spaces

$$\mu : \mathbb{R}[x^{S_1}]_{\leq d_1} \times \ldots \times \mathbb{R}[x^{S_L}]_{\leq d_L} \longrightarrow \mathbb{R}[x]_{\leq d}, \text{ where } d := \sum_i d_i S_i$$
$$(P_1, \ldots, P_L) \longmapsto P_L \cdots P_1,$$

**Theorem:** The function space $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}} = \mathrm{im}(\mu)$ is a <span style="color:yellow">semi-algebraic</span>, <span style="color:yellow">Euclidean-closed</span> subset of $\mathbb{R}[x]_{\leq d}$ of dimension $d_1 + \ldots + d_L + 1$.



$\mu : \mathbb{R}[x]_{\leq 2} \times \mathbb{R}[x^2]_{\leq 1} \to \mathbb{R}[x]_{\leq 4}$    $\mu : \mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x^2]_{\leq 1} \to \mathbb{R}[x]_{\leq 4}$
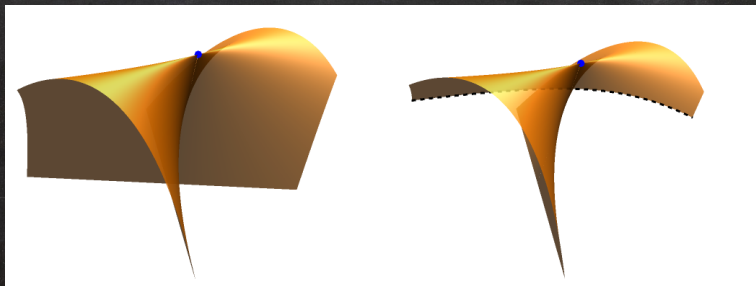
# LCN function spaces

$$\mu : \mathbb{R}[x^{S_1}]_{\leq d_1} \times \ldots \times \mathbb{R}[x^{S_L}]_{\leq d_L} \longrightarrow \mathbb{R}[x]_{\leq d}, \text{where } d := \sum_i d_i S_i$$
$$(P_1, \ldots, P_L) \longmapsto P_L \cdots P_1,$$

**Theorem:** The function space $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}} = \mathrm{im}(\mu)$ is a semi-algebraic, Euclidean-closed subset of $\mathbb{R}[x]_{\leq d}$ of dimension $d_1 + \ldots + d_L + 1$.

**Corollary:** $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}}$ is full-dimensional in $\mathbb{R}[x]_{\leq d}$ if and only if all strides $s_i = 1$.

# comparison

|  | linear<br>dense | LCN<br>$\forall i : s_i = 1$ | LCN<br>$\exists i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic &<br>full-dimensional | Euclidean closed<br>low-dimensional |

# comparison

|  | linear dense | LCN $\forall i : s_i = 1$ | LCN $\exists i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
|  |  | full-dimensional | low-dimensional |
| $\partial\mathcal{M}$ |  | non-empty | |

## comparison

|  | linear dense | LCN $\forall i : s_i = 1$ | LCN $\exists i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
| | | full-dimensional | low-dimensional |
| $\partial\mathcal{M}$ | | non-empty | |



training a network = minimizing the loss $\mathcal{L}_\mathcal{D} : \mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_\mathcal{D}} \mathbb{R}$.

# comparison

|  | linear dense | LCN $\forall i : s_i = 1$ | LCN $\exists i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
|  |  | full-dimensional | low-dimensional |
| $\partial\mathcal{M}$ |  | non-empty | |
| $\mu(\mathrm{Crit}(\mathcal{L}_{\mathcal{D}}))$ |  | often in $\partial\mathcal{M}$ | |



training a network $=$ minimizing the loss $\mathcal{L}_{\mathcal{D}} : \mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}$.

# comparison

|  | linear dense | LCN $\forall i : s_i = 1$ | LCN $\exists i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
|  | | full-dimensional | low-dimensional |
| $\partial\mathcal{M}$ | $\emptyset$ | non-empty | |
| $\mu(\mathrm{Crit}(\mathcal{L}_{\mathcal{D}}))$ | | often in $\partial\mathcal{M}$ | |



training a network $=$ minimizing the loss $\mathcal{L}_{\mathcal{D}} : \mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}$.

# comparison

| | linear dense | LCN $\forall i : s_i = 1$ | LCN $\exists i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
| | | full-dimensional | low-dimensional |
| $\partial\mathcal{M}$ | $\emptyset$ | non-empty | |
| $\mathrm{Sing}(\mathcal{M}^\circ)$ | non-empty | | |
| $\mu(\mathrm{Crit}(\mathcal{L}_\mathcal{D}))$ | often in $\mathrm{Sing}(\mathcal{M})$ | often in $\partial\mathcal{M}$ | |



training a network = minimizing the loss $\mathcal{L}_\mathcal{D} : \mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_\mathcal{D}} \mathbb{R}$.
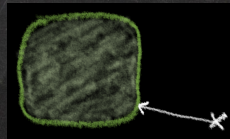
# comparison

|  | linear dense | LCN $\forall i : s_i = 1$ | LCN $\exists i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
|  |  | full-dimensional | low-dimensional |
| $\partial\mathcal{M}$ | $\emptyset$ | non-empty | |
| $\mathrm{Sing}(\mathcal{M}^\circ)$ | non-empty | $\emptyset$ | |
| $\mu(\mathrm{Crit}(\mathcal{L}_{\mathcal{D}}))$ | often in $\mathrm{Sing}(\mathcal{M})$ | often in $\partial\mathcal{M}$ | |



training a network $=$ minimizing the loss $\mathcal{L}_{\mathcal{D}} : \mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}$.
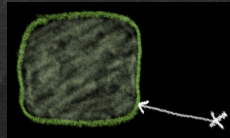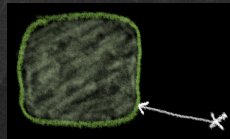
## comparison

|  | linear dense | LCN $\forall i : s_i = 1$ | LCN $\exists i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
|  |  | full-dimensional | low-dimensional |
| $\partial \mathcal{M}$ | $\emptyset$ | non-empty | non-empty |
| $\mathrm{Sing}(\mathcal{M}^\circ)$ | non-empty | $\emptyset$ | non-empty |
| $\mu(\mathrm{Crit}(\mathcal{L}_\mathcal{D}))$ | often in $\mathrm{Sing}(\mathcal{M})$ | often in $\partial \mathcal{M}$ | **??** |



training a network = minimizing the loss $\mathcal{L}_\mathcal{D} : \mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_\mathcal{D}} \mathbb{R}$.

# training with the squared error loss

Given training data $\mathcal{D} = \{(X_i, Y_i) \in \mathbb{R}^{k_0} \times \mathbb{R}^{k_L} \mid i = 1, \ldots, N\}$, the squared error loss on the function space is

$$\ell_{\mathcal{D}} : \mathbb{R}^{k_L \times k_0} \longrightarrow \mathbb{R},$$

$$T \longmapsto \sum_{i=1}^{N} \| Y_i - T X_i \|^2.$$

# training with the squared error loss

Given training data $\mathcal{D} = \{(X_i, Y_i) \in \mathbb{R}^{k_0} \times \mathbb{R}^{k_L} \mid i = 1, \ldots, N\}$, the squared error loss on the function space is

$$\ell_{\mathcal{D}} : \mathbb{R}^{k_L \times k_0} \longrightarrow \mathbb{R},$$

$$T \longmapsto \sum_{i=1}^{N} \| Y_i - TX_i \|^2.$$

Training an LCN minimizes the squared error loss on the parameter space:

$$\mathcal{L}_{\mathcal{D}} : \mathbb{R}^{d_1} \times \ldots \times \mathbb{R}^{d_L} \xrightarrow{\mu} \mathcal{M}_{\boldsymbol{d}, \boldsymbol{S}} \subseteq \mathbb{R}^{k_L \times k_0} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R},$$

$$(w_1, \ldots, w_L) \longmapsto T_{w_L, s_L} \cdots T_{w_1, s_1} \longmapsto \ell_{\mathcal{D}}(T_{w_L, s_L} \cdots T_{w_1, s_1})$$

# training LCNs with the squared error loss

$$\mathcal{L}_{\mathcal{D}} : \mathbb{R}^{d_1} \times \ldots \times \mathbb{R}^{d_L} \xrightarrow{\mu} \mathcal{M}_{\boldsymbol{d}, \boldsymbol{S}} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}$$

**Theorem**
Consider an LCN with all strides $> 1$. Let $N \geq \sum_i d_i S_i + 1$.
For almost all data $\mathcal{D} \in (\mathbb{R}^{k_0} \times \mathbb{R}^{k_L})^N$, every critical point $\boldsymbol{w}$ of $\mathcal{L}_{\mathcal{D}}$ satisfies one of the following:

# training LCNs with the squared error loss

$$\mathcal{L}_{\mathcal{D}} : \mathbb{R}^{d_1} \times \ldots \times \mathbb{R}^{d_L} \xrightarrow{\mu} \mathcal{M}_{\boldsymbol{d},\boldsymbol{S}} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}$$

**Theorem**
Consider an LCN with all strides $> 1$. Let $N \geq \sum_i d_i S_i + 1$.
For almost all data $\mathcal{D} \in (\mathbb{R}^{k_0} \times \mathbb{R}^{k_L})^N$, every critical point $\boldsymbol{w}$ of $\mathcal{L}_{\mathcal{D}}$ satisfies one of the following:

◆ $\mu(\boldsymbol{w}) = 0$, or

◆ $\mu(\boldsymbol{w})$ is a smooth, interior point of $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}}$.

# comparison

|  | linear dense | LCN $\forall i : s_i = 1$ | LCN $\forall i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
| | | full-dimensional | low-dimensional |
| $\partial\mathcal{M}$ | $\emptyset$ | non-empty | non-empty |
| $\mathrm{Sing}(\mathcal{M}^\circ)$ | non-empty | $\emptyset$ | non-empty |
| $\mu(\mathrm{Crit}(\mathcal{L}_\mathcal{D}))$ | often in $\mathrm{Sing}(\mathcal{M})$ | often in $\partial\mathcal{M}$ | almost never in $\mathrm{Sing}(\mathcal{M}^\circ)$ or $\partial\mathcal{M}$ |

training a network $=$ minimizing the loss $\mathcal{L}_\mathcal{D} : \mathbb{R}^N \xrightarrow{\;\mu\;} \mathcal{M} \xrightarrow{\;\ell_\mathcal{D}\;} \mathbb{R}$.

# comparison

|  | linear dense | LCN $\forall i : s_i = 1$ | LCN $\forall i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
|  |  | full-dimensional | low-dimensional |
| $\partial\mathcal{M}$ | $\emptyset$ | non-empty | non-empty |
| $\mathrm{Sing}(\mathcal{M}^\circ)$ | non-empty | $\emptyset$ | non-empty |
| $\mu(\mathrm{Crit}(\mathcal{L}_\mathcal{D}))$ | often in $\mathrm{Sing}(\mathcal{M})$ | often in $\partial\mathcal{M}$ | almost never in $\mathrm{Sing}(\mathcal{M}^\circ)$ or $\partial\mathcal{M}$ |

training a network $=$ minimizing the loss $\mathcal{L}_\mathcal{D} : \mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_\mathcal{D}} \mathbb{R}$.

A critical point $\theta \in \mathrm{Crit}(\mathcal{L}_\mathcal{D})$ is called **spurious** if $\mu(\theta) \notin \mathrm{Crit}(\ell_\mathcal{D})$.

## comparison

|  | linear dense | LCN $\forall i : s_i = 1$ | LCN $\forall i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
| | | full-dimensional | low-dimensional |
| $\partial\mathcal{M}$ | $\emptyset$ | non-empty | non-empty |
| $\mathrm{Sing}(\mathcal{M}^\circ)$ | non-empty | $\emptyset$ | non-empty |
| $\mu(\mathrm{Crit}(\mathcal{L}_{\mathcal{D}}))$ | often in $\mathrm{Sing}(\mathcal{M})$ | often in $\partial\mathcal{M}$ | almost never in $\mathrm{Sing}(\mathcal{M}^\circ)$ or $\partial\mathcal{M}$ |
| critical points spurious? | often | | |

training a network = minimizing the loss $\mathcal{L}_{\mathcal{D}} : \mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}$.

A critical point $\theta \in \mathrm{Crit}(\mathcal{L}_{\mathcal{D}})$ is called **spurious** if $\mu(\theta) \notin \mathrm{Crit}(\ell_{\mathcal{D}})$.

# comparison

| | linear dense | LCN $\forall i : s_i = 1$ | LCN $\forall i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
| | | full-dimensional | low-dimensional |
| $\partial \mathcal{M}$ | $\emptyset$ | non-empty | non-empty |
| $\mathrm{Sing}(\mathcal{M}^\circ)$ | non-empty | $\emptyset$ | non-empty |
| $\mu(\mathrm{Crit}(\mathcal{L}_\mathcal{D}))$ | often in $\mathrm{Sing}(\mathcal{M})$ | often in $\partial \mathcal{M}$ | almost never in $\mathrm{Sing}(\mathcal{M}^\circ)$ or $\partial \mathcal{M}$ |
| critical points spurious? | often | often | almost never |

training a network = minimizing the loss $\mathcal{L}_\mathcal{D} : \mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_\mathcal{D}} \mathbb{R}$.

A critical point $\theta \in \mathrm{Crit}(\mathcal{L}_\mathcal{D})$ is called **spurious** if $\mu(\theta) \notin \mathrm{Crit}(\ell_\mathcal{D})$.

# training LCNs with the squared error loss

$$\mathcal{L}_{\mathcal{D}} : \mathbb{R}^{d_1} \times \ldots \times \mathbb{R}^{d_L} \xrightarrow{\mu} \mathcal{M}_{\boldsymbol{d,S}} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}$$

**Theorem**
Consider an LCN with all strides $> 1$. Let $N \geq \sum_i d_i S_i + 1$.
For almost all data $\mathcal{D} \in (\mathbb{R}^{k_0} \times \mathbb{R}^{k_L})^N$, every critical point $\boldsymbol{w}$ of $\mathcal{L}_{\mathcal{D}}$ satisfies one of the following:

# training LCNs with the squared error loss

$$\mathcal{L}_{\mathcal{D}} : \mathbb{R}^{d_1} \times \ldots \times \mathbb{R}^{d_L} \xrightarrow{\mu} \mathcal{M}_{\boldsymbol{d},\boldsymbol{S}} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}$$

**Theorem**
Consider an LCN with all strides $> 1$. Let $N \geq \sum_i d_i S_i + 1$.
For almost all data $\mathcal{D} \in (\mathbb{R}^{k_0} \times \mathbb{R}^{k_L})^N$, every critical point $\boldsymbol{w}$ of $\mathcal{L}_{\mathcal{D}}$ satisfies one of the following:

- $\mu(\boldsymbol{w}) = 0$, or
- $\mu(\boldsymbol{w})$ is a smooth, interior point of $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}}$ and $\boldsymbol{w}$ is a regular point of $\mu$.

# training LCNs with the squared error loss

$$\mathcal{L}_{\mathcal{D}} : \mathbb{R}^{d_1} \times \ldots \times \mathbb{R}^{d_L} \xrightarrow{\mu} \mathcal{M}_{\boldsymbol{d}, \boldsymbol{S}} \xrightarrow{\ell_{\mathcal{D}}} \mathbb{R}$$

**Theorem**
Consider an LCN with all strides $> 1$. Let $N \geq \sum_i d_i S_i + 1$.
For almost all data $\mathcal{D} \in (\mathbb{R}^{k_0} \times \mathbb{R}^{k_L})^N$, every critical point $\boldsymbol{w}$ of $\mathcal{L}_{\mathcal{D}}$ satisfies one of the following:

- ◆ $\mu(\boldsymbol{w}) = 0$, or
- ◆ $\mu(\boldsymbol{w})$ is a smooth, interior point of $\mathcal{M}_{\boldsymbol{d}, \boldsymbol{S}}$ and $\boldsymbol{w}$ is a regular point of $\mu$. In particular, $\mu(\boldsymbol{w})$ is a critical point of $\ell_{\mathcal{D}}|_{\mathrm{Reg}(\mathcal{M}_{\boldsymbol{d}, \boldsymbol{S}}^{\circ})}$.

# reducing LCNs

$\mu : \mathbb{R}[x]_{\leq 2} \times \mathbb{R}[x^2]_{\leq 1} \to \mathbb{R}[x]_{\leq 4}$     $\mu : \mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x^2]_{\leq 1} \to \mathbb{R}[x]_{\leq 4}$



$$\mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x]_{\leq 1} \qquad \times \qquad \mathbb{R}[x^2]_{\leq 1} \longrightarrow \mathcal{M}_{(1,1,1),(1,1,2)}$$

$$\downarrow \qquad\qquad\qquad\qquad\qquad\qquad \downarrow \qquad\qquad\qquad \downarrow$$
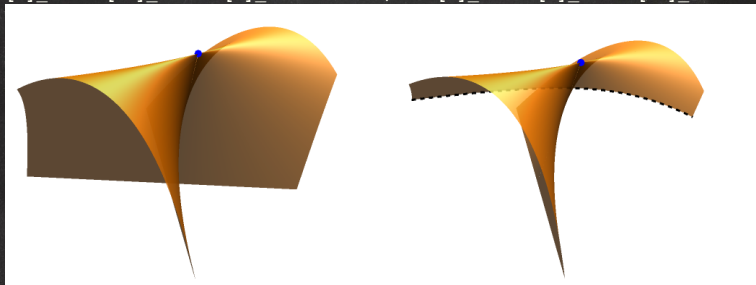
$$\mathbb{R}[x]_{\leq 2} \qquad\qquad \times \qquad\qquad \mathbb{R}[x^2]_{\leq 1} \longrightarrow \mathcal{M}_{(2,1),(1,2)}$$

# reducing LCNs

$\mu : \mathbb{R}[x]_{\leq 2} \times \mathbb{R}[x^2]_{\leq 1} \to \mathbb{R}[x]_{\leq 4}$ $\qquad\qquad$ $\mu : \mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x^2]_{\leq 1} \to \mathbb{R}[x]_{\leq 4}$



$$
\begin{array}{ccccc}
\mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x]_{\leq 1} & \times & \mathbb{R}[x^2]_{\leq 1} & \longrightarrow & \mathcal{M}_{(1,1,1),(1,1,2)} \\
\downarrow & & \downarrow & & \downarrow \\
\mathbb{R}[x]_{\leq 2} & \times & \mathbb{R}[x^2]_{\leq 1} & \longrightarrow & \mathcal{M}_{(2,1),(1,2)}
\end{array}
$$

Given an LCN $(\boldsymbol{d}, \boldsymbol{S})$, merging neighboring layers with the same $S_i$ yields an LCN $(\tilde{\boldsymbol{d}}, \tilde{\boldsymbol{S}})$ with $1 = \tilde{S}_1 < \tilde{S}_2 < \tilde{S}_3 < \ldots$ (i.e., all strides $> 1$), called the reduced LCN.

# Singularities

**Lemma:** $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}} \subseteq \mathcal{M}_{\tilde{\boldsymbol{d}},\tilde{\boldsymbol{S}}}$ and $\overline{\mathcal{M}}_{\boldsymbol{d},\boldsymbol{S}} = \overline{\mathcal{M}}_{\tilde{\boldsymbol{d}},\tilde{\boldsymbol{S}}}$,
where $\overline{\phantom{x}}$ denotes the Zariski closure inside $\mathbb{R}[x]_{\leq d}$.

# Singularities

**Lemma:** $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}} \subseteq \mathcal{M}_{\tilde{\boldsymbol{d}},\tilde{\boldsymbol{S}}}$ and $\overline{\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}}} = \overline{\mathcal{M}}_{\tilde{\boldsymbol{d}},\tilde{\boldsymbol{S}}}$,
where $\overline{\phantom{.}}$ denotes the Zariski closure inside $\mathbb{R}[x]_{\leq d}$.

**Theorem** Let $(\boldsymbol{d}, \boldsymbol{S})$ be a reduced LCN with $L$ layers.

- If $L = 1$ (i.e., any associated non-reduced LCN has all strides equal 1),
  then $\overline{\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}}} = \mathbb{R}[x]_{\leq d}$.

# Singularities

**Lemma:** $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}} \subseteq \mathcal{M}_{\tilde{\boldsymbol{d}},\tilde{\boldsymbol{S}}}$ and $\overline{\mathcal{M}}_{\boldsymbol{d},\boldsymbol{S}} = \overline{\mathcal{M}}_{\tilde{\boldsymbol{d}},\tilde{\boldsymbol{S}}}$,
where $\bar{\ }$ denotes the Zariski closure inside $\mathbb{R}[x]_{\leq d}$.

**Theorem** Let $(\boldsymbol{d}, \boldsymbol{S})$ be a reduced LCN with $L$ layers.

- If $L = 1$ (i.e., any associated non-reduced LCN has all strides equal 1), then $\overline{\mathcal{M}}_{\boldsymbol{d},\boldsymbol{S}} = \mathbb{R}[x]_{\leq d}$.

- If $L > 1$, $\deg \overline{\mathcal{M}}_{\boldsymbol{d},\boldsymbol{S}} > 1$ and

$$\mathrm{Sing}(\overline{\mathcal{M}}_{\boldsymbol{d},\boldsymbol{S}}) = \{0\} \cup \bigcup_{\boldsymbol{d}' \in D} \overline{\mathcal{M}}_{\boldsymbol{d}',\boldsymbol{S}} = \{0\} \cup \bigcup_{\boldsymbol{d}' \in D} \mathcal{M}_{\boldsymbol{d}',\boldsymbol{S}},$$

where $D := \{\boldsymbol{d}' \in \mathbb{Z}_{\geq 0}^{L} \mid \overline{\mathcal{M}}_{\boldsymbol{d}',\boldsymbol{S}} \subsetneq \overline{\mathcal{M}}_{\boldsymbol{d},\boldsymbol{S}}\}$

# Singularities

**Lemma:** $\mathcal{M}_{\boldsymbol{d}, \boldsymbol{S}} \subseteq \mathcal{M}_{\tilde{\boldsymbol{d}}, \tilde{\boldsymbol{S}}}$ and $\overline{\mathcal{M}}_{\boldsymbol{d}, \boldsymbol{S}} = \overline{\mathcal{M}}_{\tilde{\boldsymbol{d}}, \tilde{\boldsymbol{S}}}$,
where $\bar{\cdot}$ denotes the Zariski closure inside $\mathbb{R}[x]_{\leq d}$.

**Theorem** Let $(\boldsymbol{d}, \boldsymbol{S})$ be a reduced LCN with $L$ layers.

- If $L = 1$ (i.e., any associated non-reduced LCN has all strides equal 1), then $\overline{\mathcal{M}}_{\boldsymbol{d}, \boldsymbol{S}} = \mathbb{R}[x]_{\leq d}$.
- If $L > 1$, $\deg \overline{\mathcal{M}}_{\boldsymbol{d}, \boldsymbol{S}} > 1$ and

$$\mathrm{Sing}(\overline{\mathcal{M}}_{\boldsymbol{d}, \boldsymbol{S}}) = \{0\} \cup \bigcup_{\boldsymbol{d}' \in D} \overline{\mathcal{M}}_{\boldsymbol{d}', \boldsymbol{S}} = \{0\} \cup \bigcup_{\boldsymbol{d}' \in D} \mathcal{M}_{\boldsymbol{d}', \boldsymbol{S}},$$

where $D := \{\boldsymbol{d}' \in \mathbb{Z}_{\geq 0}^L \mid \overline{\mathcal{M}}_{\boldsymbol{d}', \boldsymbol{S}} \subsetneq \overline{\mathcal{M}}_{\boldsymbol{d}, \boldsymbol{S}}\}$
$= \{\boldsymbol{d}' \in \mathbb{Z}_{\geq 0}^L \mid \boldsymbol{d}' \neq \boldsymbol{d}, \sum_{i=1}^L d_i' S_i = \sum_{i=1}^L d_i S_i, \forall l : \sum_{i=l}^L d_i' S_i \geq \sum_{i=l}^L d_i S_i\}$

# Example

$\mu : \mathbb{R}[x]_{\leq 2} \times \mathbb{R}[x^2]_{\leq 1} \to \mathbb{R}[x]_{\leq 4}$ $\qquad$ $\mu : \mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x^2]_{\leq 1} \to \mathbb{R}[x]_{\leq 4}$



$\mathbb{R}[x]_{\leq 2} \times \mathbb{R}[x^2]_{\leq 1} \to \mathcal{M}_{(2,1),(1,2)}$

$\mathrm{Sing}(\overline{\mathcal{M}_{(2,1),(1,2)}}) =$

# Example

$\mu : \mathbb{R}[x]_{\leq 2} \times \mathbb{R}[x^2]_{\leq 1} \to \mathbb{R}[x]_{\leq 4}$ $\qquad$ $\mu : \mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x]_{\leq 1} \times \mathbb{R}[x^2]_{\leq 1} \to \mathbb{R}[x]_{\leq 4}$
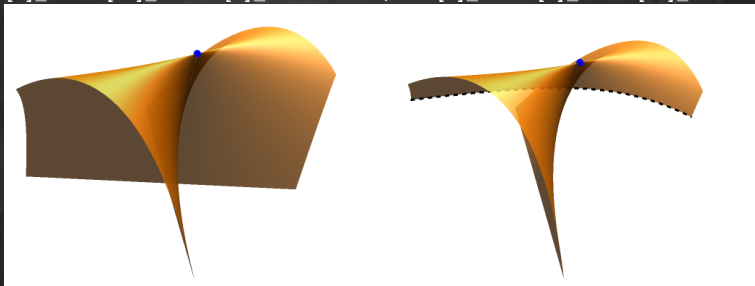


$\mathbb{R}[x]_{\leq 2} \times \mathbb{R}[x^2]_{\leq 1} \to \mathcal{M}_{(2,1),(1,2)}$
$\mathrm{Sing}(\overline{\mathcal{M}_{(2,1),(1,2)}}) = \mathcal{M}_{(0,2),(1,2)} = \mathbb{R}[x^2]_{\leq 2}$

# Relative Boundary

$\partial \mathcal{M}_{d,S} =$ points in $\mathcal{M}_{d,S}$ that are limits of sequences in $\overline{\mathcal{M}}_{d,S} \setminus \mathcal{M}_{d,S}$.

# Relative Boundary

$\partial \mathcal{M}_{d,S}$ = points in $\mathcal{M}_{d,S}$ that are limits of sequences in $\overline{\mathcal{M}}_{d,S} \setminus \mathcal{M}_{d,S}$.

**Recall:** $\mathcal{M}_{d,S} \subseteq \mathcal{M}_{\tilde{d},\tilde{S}} \subseteq \overline{\mathcal{M}}_{d,S} = \overline{\mathcal{M}}_{\tilde{d},\tilde{S}}$

# Relative Boundary

$\partial \mathcal{M}_{d,S}$ = points in $\mathcal{M}_{d,S}$ that are limits of sequences in $\overline{\mathcal{M}}_{d,S} \setminus \mathcal{M}_{d,S}$.

**Recall:** $\mathcal{M}_{d,S} \subseteq \mathcal{M}_{\tilde{d},\tilde{S}} \subseteq \overline{\mathcal{M}}_{d,S} = \overline{\mathcal{M}}_{\tilde{d},\tilde{S}}$
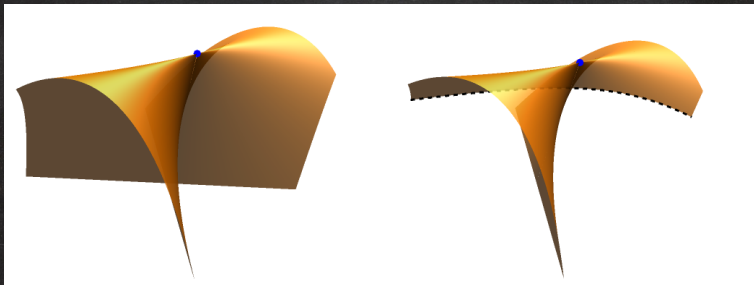
- ◆ reduced boundary points: limits in $\mathcal{M}_{d,S}$ of sequences in $\overline{\mathcal{M}}_{d,S} \setminus \mathcal{M}_{\tilde{d},\tilde{S}}$
- ◆ stride-1 boundary points: limits in $\mathcal{M}_{d,S}$ of sequences in $\mathcal{M}_{\tilde{d},\tilde{S}} \setminus \mathcal{M}_{d,S}$

# Relative Boundary

$\partial \mathcal{M}_{\boldsymbol{d},\boldsymbol{S}} =$ points in $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}}$ that are limits of sequences in $\overline{\mathcal{M}}_{\boldsymbol{d},\boldsymbol{S}} \setminus \mathcal{M}_{\boldsymbol{d},\boldsymbol{S}}$.

**Recall:** $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}} \subseteq \mathcal{M}_{\tilde{\boldsymbol{d}},\tilde{\boldsymbol{S}}} \subseteq \overline{\mathcal{M}}_{\boldsymbol{d},\boldsymbol{S}} = \overline{\mathcal{M}}_{\tilde{\boldsymbol{d}},\tilde{\boldsymbol{S}}}$

- ◆ reduced boundary points: limits in $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}}$ of sequences in $\overline{\mathcal{M}}_{\boldsymbol{d},\boldsymbol{S}} \setminus \mathcal{M}_{\tilde{\boldsymbol{d}},\tilde{\boldsymbol{S}}}$
- ◆ stride-1 boundary points: limits in $\mathcal{M}_{\boldsymbol{d},\boldsymbol{S}}$ of sequences in $\mathcal{M}_{\tilde{\boldsymbol{d}},\tilde{\boldsymbol{S}}} \setminus \mathcal{M}_{\boldsymbol{d},\boldsymbol{S}}$



reduced boundary points have at least codimension 2
stride-1 boundary points (if existent) have codimension 1

# comparison

| | linear dense | LCN $\forall i : s_i = 1$ | LCN $\forall i : s_i > 1$ |
|---|---|---|---|
| $\mathcal{M}$ | algebraic variety | semialgebraic & Euclidean closed | |
| | | full-dimensional | low-dimensional |
| $\partial\mathcal{M}$ | $\emptyset$ | non-empty | non-empty |
| $\mathrm{Sing}(\mathcal{M}^\circ)$ | non-empty | $\emptyset$ | non-empty |
| $\mu(\mathrm{Crit}(\mathcal{L_D}))$ | often in $\mathrm{Sing}(\mathcal{M})$ | often in $\partial\mathcal{M}$ | almost never in $\mathrm{Sing}(\mathcal{M}^\circ)$ or $\partial\mathcal{M}$ |
| critical points spurious? | often | often | almost never |

training a network = minimizing the loss $\mathcal{L_D} : \mathbb{R}^N \xrightarrow{\mu} \mathcal{M} \xrightarrow{\ell_D} \mathbb{R}$.

A critical point $\theta \in \mathrm{Crit}(\mathcal{L_D})$ is called **spurious** if $\mu(\theta) \notin \mathrm{Crit}(\ell_D)$.