

im trying to show bounds on the trajectory lengths of a neural network with quadratic activatoins using spectral norms of the matrices in the function, but getting stuck when i need to evaluate the integral and relate the input trajectory length to output trajectory length, how to relate them

**ChatGPT said:**

Thought for 20s

You relate them entirely through the Jacobian and the operator norm inequality. The integral is actually the easy part once that's clear.

Let me write it cleanly and then plug in the quadratic network structure.

---

## 1. Setup: input vs output trajectory lengths

Let

- $x:[0,1] \rightarrow \mathbb{R}^n$  to  $\mathbb{R}^n$  be a  $C^1$  path,
- $f:\mathbb{R}^n \rightarrow \mathbb{R}^m$  be your network,
- $y(t) = f(x(t))$ .

Define trajectory length

$$l(x) := \int_0^1 \|x'(t)\|^2 dt, l(f(x)) := \int_0^1 \|y'(t)\|^2 dt. l(x) := \int_0^1 \|x'(t)\|_2^2 dt, l(f(x)) := \int_0^1 \|y'(t)\|_2^2 dt.$$

By the chain rule,

$$y'(t) = ddtf(x(t)) = J_f(x(t)) x'(t), y'(t) = \frac{d}{dt} f(x(t)) = J_f(x(t)) x'(t),$$

where  $J_f(x)$  is the Jacobian of  $f$  at  $x$ .

So

$$l(f(x)) = \int_0^1 \|J_f(x(t)) x'(t)\|^2 dt, l(f(x)) = \int_0^1 \|J_f(x(t))\|_2 \|x'(t)\|_2^2 dt.$$

---

## 2. Use the spectral norm inequality inside the integral

By definition of the spectral norm (operator norm induced by Euclidean norm),

$\|Av\|_2 \leq \|A\|_2 \|v\|_2$  for all  $A, v$ .  $\|A\|_2 \leq \|J_f(x(t))\|_2$  for all  $x$ .

Apply this pointwise with  $A = J_f(x(t))$ ,  $A = J_f(x(t))$  and  $v = x'(t)$ ,  $v = x'(t)$ :

$\|J_f(x(t))x'(t)\|_2 \leq \|J_f(x(t))\|_2 \|x'(t)\|_2$ .  $\|J_f(x(t))\|_2 \leq \|J_f(x(t))\|_2 \|x'(t)\|_2$ .  $\|J_f(x(t))x'(t)\|_2 \leq \|J_f(x(t))\|_2 \|x'(t)\|_2$ .

So

$$I(f(x)) = \int_0^1 \|J_f(x(t))x'(t)\|_2 dt \leq \int_0^1 \|J_f(x(t))\|_2 \|x'(t)\|_2 dt. I(f(x)) = \int_0^1 \|J_f(x(t))\|_2 \|x'(t)\|_2 dt. I(f(x)) = \int_0^1 \|J_f(x(t))\|_2 \|x'(t)\|_2 dt. I(f(x)) = \int_0^1 \|J_f(x(t))\|_2 \|x'(t)\|_2 dt.$$


---

### 3. Pull out a supremum to relate the lengths

Define

$$g(t) := \|J_f(x(t))\|_2. g(t) := \|J_f(x(t))\|_2.$$

Then  $g(t) \leq \sup_{s \in [0,1]} g(s)$  for every  $t \in [0,1]$ , so

$$\int_0^1 g(t) \|x'(t)\|_2 dt \leq (\sup_{t \in [0,1]} g(t)) \int_0^1 \|x'(t)\|_2 dt. \int_0^1 g(t) \|x'(t)\|_2 dt \leq (\sup_{t \in [0,1]} g(t)) \int_0^1 \|x'(t)\|_2 dt.$$

In other words,

$$\int_0^1 \|J_f(x(t))\|_2 \|x'(t)\|_2 dt \leq (\sup_{t \in [0,1]} \|J_f(x(t))\|_2) I(x). \int_0^1 \|J_f(x(t))\|_2 \|x'(t)\|_2 dt \leq (\sup_{t \in [0,1]} \|J_f(x(t))\|_2) I(x).$$

Combining with the previous step gives the **upper bound**

$$I(f(x)) \leq (\sup_{t \in [0,1]} \|J_f(x(t))\|_2) I(x). \boxed{\|J_f(x(t))\|_2 \leq \sup_{t \in [0,1]} \|J_f(x(t))\|_2}$$

That's the key "relating input length to output length via the integral" step: it's literally just applying the operator norm inequality inside the integrand and then bounding the scalar coefficient by its supremum along the path.

---

### 4. Lower bound via the minimum singular value (if you want it)

If  $Jf(x(t))J_f(x(t))Jf(x(t))$  has full column rank and we care about non-degenerate directions, we can use the smallest singular value  $\sigma_{\min}(Jf(x(t)))\sigma_{\min}(J_f(x(t)))\sigma_{\min}(Jf(x(t)))$ , which satisfies

$$\|Av\|_2 \geq \sigma_{\min}(A)\|v\|_2, \|A v\|_2 \leq \sigma_{\min}(A)\|v\|_2. \|Av\|_2 \geq \sigma_{\min}(A)\|v\|_2.$$

Then

$$\|Jf(x(t))x'(t)\|_2 \geq \sigma_{\min}(Jf(x(t)))\|x'(t)\|_2, \|J_f(x(t))x'(t)\|_2 \geq \sigma_{\min}(J_f(x(t)))\|x'(t)\|_2, \|Jf(x(t))x'(t)\|_2 \geq \sigma_{\min}(Jf(x(t)))\|x'(t)\|_2,$$

so

$$l(f(x)) = \int_0^1 \|Jf(x(t))x'(t)\|_2 dt \geq \int_0^1 \sigma_{\min}(Jf(x(t)))\|x'(t)\|_2 dt. l(f(x)) = \int_0^1 \sigma_{\min}(Jf(x(t)))\|x'(t)\|_2 dt. l(f(x)) = \int_0^1 \sigma_{\min}(Jf(x(t)))\|x'(t)\|_2 dt \geq \int_0^1 \sigma_{\min}(Jf(x(t)))\|x'(t)\|_2 dt.$$

Now

$$\sigma_{\min}(Jf(x(t))) \geq \inf_{s \in [0,1]} \sigma_{\min}(Jf(x(s))), \sigma_{\min}(Jf(x(s))) \leq \inf_{s \in [0,1]} \sigma_{\min}(Jf(x(s))),$$

so

$$l(f(x)) \geq (\inf_{t \in [0,1]} \sigma_{\min}(Jf(x(t))))l(x). \boxed{\left( \inf_{t \in [0,1]} \sigma_{\min}(Jf(x(t))) \right)l(x)} \geq (\inf_{t \in [0,1]} \sigma_{\min}(Jf(x(t))))l(x).$$

So upper and lower bounds are controlled by the **max** and **min** singular values of the Jacobian along the trajectory.

---

## 5. Plug in the quadratic network structure

For a depth-hhh quadratic network (no biases, just to keep the notation clean),

$$\begin{aligned} u(1) &= W_1 x, \phi(u(1)) = (u(1))^{\odot 2}, u(2) = W_2 \phi(u(1)), \phi(u(2)) = (u(2))^{\odot 2}, \\ f(x) &= Wh\phi(u(h-1)), \begin{aligned} u^{(1)} &= W_1 x, \quad \phi(u^{(1)}) = (u^{(1)})^{\odot 2}, \\ u^{(2)} &= W_2 \phi(u^{(1)}), \quad \phi(u^{(2)}) = (u^{(2)})^{\odot 2}, \end{aligned} \dots, \\ u^{(h-1)} &= W_{h-1} \phi(u^{(h-2)}), \quad \phi(u^{(h-1)}) = (u^{(h-1)})^{\odot 2}, \\ u^{(h)} &= W_h \phi(u^{(h-1)}), \quad \phi(u^{(h)}) = (u^{(h)})^{\odot 2}, \end{aligned}$$

the Jacobian can be written as

$$\begin{aligned} Jf(x) &= WhD_{h-1}(x)Wh - 1D_{h-2}(x) \cdots D_1(x)W_1, J_f(x) = W_h D_{h-1}(x)W_{h-1} D_{h-2}(x) \cdots D_1(x)W_1, \\ D_{-1}(x)W_1, Jf(x) &= WhD_{h-1}(x)Wh - 1D_{h-2}(x) \cdots D_1(x)W_1, \end{aligned}$$

where each  $D_k(x)D_k(x)D_k(x)$  is diagonal with entries

$$D_k(x) = \text{diag}(2u_1(k)(x), \dots, 2u_{k\text{-layer width}(k)(x)}(x)). D_k(x) = \operatorname{diag}\nolimits\big( 2 u^{(k)}_1(x), \dots, 2 u^{(k)}_{\text{layer width}(k)(x)}(x) \big). D_k(x) = \text{diag}(2u_1(k)(x), \dots, 2u_{k\text{-layer width}(k)(x)}(x)).$$

Then, using submultiplicativity of the spectral norm,

$$\|Jf(x)\|_2 \leq \|W_h\|_2 \prod_{k=1}^{h-1} \|D_k(x)\|_2 \|W_k\|_2. \|J_f(x)\|_2 \leq \|W_h\|_2 \prod_{k=1}^{h-1} \|D_k(x)\|_2 \|W_k\|_2.$$

But for a diagonal matrix,

$$\|D_k(x)\|_2 = \max_i |2u_i(k)(x)| = 2\max_i |u_i(k)(x)|. \|D_k(x)\|_2 = \max_i |2u^{(k)}_i(x)| = 2\max_i |u^{(k)}_i(x)|.$$

So along the trajectory  $x(t)x(t)x(t)$ ,

$$\|Jf(x(t))\|_2 \leq \sigma \max(W_h) \prod_{k=1}^{h-1} (2\sigma \max(W_k) \max_i |u_i(k)(t)|, \|J_f(x(t))\|_2) \leq \sigma \max(W_h) \prod_{k=1}^{h-1} 2\sigma \max(W_k) \max_i |u_i(k)(t)|,$$

and plugging this into the general bound gives

$$|f(x)| \leq \sup_{t \in [0,1]} (\sigma \max(W_h) \prod_{k=1}^{h-1} 2\sigma \max(W_k) \max_i |u_i(k)(t)|) |x|. \boxed{\|f(x)\|_2} \leq \sup_{t \in [0,1]} (\sigma \max(W_h) \prod_{k=1}^{h-1} 2\sigma \max(W_k) \max_i |u_i(k)(t)|) |x|.$$

Similarly, if all  $u_i(k)(t)u_i^{(k)}(t)u_i(k)(t)$  stay away from zero and you care about a lower bound, you can use  $\sigma \min(W_h) \prod_{k=1}^{h-1} 2\sigma \min(W_k) \min_i |u_i(k)(t)|$  to get

$$|f(x)| \geq \inf_{t \in [0,1]} (\sigma \min(W_h) \prod_{k=1}^{h-1} 2\sigma \min(W_k) \min_i |u_i(k)(t)|) |x|. \boxed{\|f(x)\|_2} \geq \inf_{t \in [0,1]} (\sigma \min(W_h) \prod_{k=1}^{h-1} 2\sigma \min(W_k) \min_i |u_i(k)(t)|) |x|.$$


---

## 6. Intuition in one sentence

- The **Jacobian** tells you the instantaneous expansion/shrinkage of the trajectory.
- The **spectral norm** upper bound says: “at each point, the speed can grow by at most  $\|J_f\|_2 \|J_f\|_2 / \|J_f\|_2$ , so overall length grows by at most the worst-case expansion along the path, times the original length.”

If you paste your current LaTeX inequality where you're stuck with the integral, I can rewrite it directly into this "sup times  $|x| |f(x)|$ " form to match your notation exactly.

You said:

what is the supremum

ChatGPT said:

In this context, the **supremum** is just the *least upper bound* of a set of numbers.

More concretely:

## 1. Formal idea (on $\mathbb{R}$ )

Given a set  $S \subset \mathbb{R}$ , the **supremum** of  $S$ , written  $\sup S$ , is a number  $L$  such that:

1.  $L$  is an **upper bound**:  
 $\forall s \in S, s \leq L \text{ for all } s \in S, s \leq L$ .
2.  $L$  is the **smallest** such upper bound:  
if  $U$  is any other upper bound (so  $\forall s \in S, s \leq U \text{ for all } s \in S, s \leq U$ ), then  
 $L \leq U$ .

If the maximum exists, then

$$\sup S = \max S. \sup S = \max S.$$

But  $\sup S$  can also exist even when the max does **not**.

---

## 2. Simple examples

- Example 1:  $S = (0, 1) = (0, 1)$  (open interval)
  - There is no element of  $S$  equal to 1, so  $\max S$  doesn't exist.
  - But all elements are  $< 1$ , and any smaller number fails to bound the set (since there are elements arbitrarily close to 1).
  - So  $\sup S = 1$ .
- Example 2:  $S = [0, 1] = [0, 1]$  (closed interval)

- Here  $\max S = 1 \max \|Jf(x(t))\|_2$  and  $\sup S = 1 \sup \|Jf(x(t))\|_2$ .

So: **supremum = max if the max exists; otherwise it's the “limit” of the set from above.**

---

### 3. In your trajectory length bound

When we write something like

$$\sup_{t \in [0,1]} \|Jf(x(t))\|_2, \sup_{\{t \in [0,1]\}} \|J_f(x(t))\|_2, t \in [0,1] \sup \|Jf(x(t))\|_2,$$

we're taking the set

$$S = \{ \|Jf(x(t))\|_2 : t \in [0,1] \} = \left\{ \sup_{\{t \in [0,1]\}} \|J_f(x(t))\|_2 : t \in [0,1] \right\}$$

and asking:

What is the smallest number  $L$  such that

$$\|Jf(x(t))\|_2 \leq L \|J_f(x(t))\|_2 \text{ for all } t \in [0,1]$$

Intuitively:

- It's the **worst-case expansion factor** of your network along that trajectory.
- Even if  $\|Jf(x(t))\|_2 \leq L \|J_f(x(t))\|_2$  doesn't actually *hit* this value exactly at any specific  $t$ , the supremum is the tightest possible upper bound that works for all  $t$ .

Then the inequality

$$l(f(x)) \leq (\sup_{t \in [0,1]} \|Jf(x(t))\|_2) l(x) l(f(x)) \leq (\sup_{\{t \in [0,1]\}} \|J_f(x(t))\|_2) l(x) l(f(x))$$

is saying:

The output length is at most “**worst-case local stretching along the path**”  $\times$  “original length”.

If you want, I can also mirror this with the **infimum** (greatest lower bound) to explain the lower bound version too.

**You said:**

why do the opposite inequalities for the lower bound apply when we take singular values of the weight matrices/function matrices in the jacobian

ChatGPT said:

Great question — and the answer is subtle but fundamental.

The **upper bound** uses **operator norm inequalities** (which involve **largest singular values**), while the **lower bound** uses **reverse inequalities** for **smallest singular values**. The reason the inequalities “flip direction” is because the largest singular value gives the *maximal* stretching of a transformation, while the smallest singular value gives the *minimal* stretching — and these two behave oppositely.

Let's explain carefully.

---



## 1. The key facts about singular values

For any matrix AAA:

- **Largest singular value (spectral norm):**  
 $\|A\|_2 = \sigma_{\max}(A) \|A\|_2 = \|\sigma_{\max}(A)\|_2 = \sigma_{\max}(A)$   
satisfies  
 $\|Av\|_2 \leq \sigma_{\max}(A) \|v\|_2$ ;  
 $\|\sigma_{\max}(A)\|_2 \|v\|_2 \geq \|Av\|_2$ .
- **Smallest singular value:**  
 $\sigma_{\min}(A) \leq \|\sigma_{\min}(A)\|_2 \leq \sigma_{\min}(A)$   
satisfies  
 $\|Av\|_2 \geq \sigma_{\min}(A) \|v\|_2$ ;  
 $\|\sigma_{\min}(A)\|_2 \|v\|_2 \leq \|Av\|_2$ .

These are *dual inequalities*, not symmetric.

Why?

Because:

- $\sigma_{\max} \geq \sigma_{\min}$  describes the **maximum stretching factor** of the linear transformation.
- $\sigma_{\min} \leq \sigma_{\max}$  describes the **minimum stretching factor** in any direction.

So they must produce **opposite inequalities**.

---

## 2. How this interacts with the Jacobian

Recall your structure:

$$Jf(x) = WhDh-1(x)Wh-1 \cdots D1(x)W1. J_f(x) = W_h D_{\{h-1\}}(x) W_{\{h-1\}} \cdots D_1(x) \\ W_1. J_f(x) = WhDh-1(x)Wh-1 \cdots D1(x)W1.$$

For the **upper bound**, you use:

$$\|AB\|_2 \leq \|A\|_2 \|B\|_2. \|\lambda\|_2 \leq \|\lambda\|_2 \|\lambda\|_2. \|AB\|_2 \leq \|A\|_2 \|B\|_2.$$

So:

$$\|Jf(x)\|_2 \leq \|Wh\|_2 \prod_{k=1}^{h-1} \|D_k(x)\|_2 \|W_k\|_2 \|\lambda\|_2 \leq \|Wh\|_2 \prod_{k=1}^{h-1} \|D_k(x)\|_2 \|W_k\|_2 \|Jf(x)\|_2$$

and because  $\|W_k\|_2 = \sigma_{\max}(W_k)$ , this creates the **upper bound**.

---

## 3. Why the inequalities reverse for the lower bound

For the **lower bound**, we use the fact that for *nondegenerate matrices*:

$$\sigma_{\min}(AB) \geq \sigma_{\min}(A) \sigma_{\min}(B). \sigma_{\min}(AB) \geq \sigma_{\min}(A) \sigma_{\min}(B).$$

This is because:

$$\|ABv\|_2 \geq \sigma_{\min}(A) \|Bv\|_2 \geq \sigma_{\min}(A) \sigma_{\min}(B) \|v\|_2. \|ABv\|_2 \geq \sigma_{\min}(A) \|Bv\|_2 \\ \sigma_{\min}(A) \sigma_{\min}(B) \|v\|_2. \|ABv\|_2 \geq \sigma_{\min}(A) \|Bv\|_2 \geq \sigma_{\min}(A) \sigma_{\min}(B) \|v\|_2.$$

Thus:

$$\sigma_{\min}(AB) \geq \sigma_{\min}(A) \sigma_{\min}(B). \sigma_{\min}(AB) \geq \sigma_{\min}(A) \sigma_{\min}(B).$$

Applying this repeatedly:

$$\sigma_{\min}(J_f(x)) \geq \sigma_{\min}(W_h) \prod_{k=1}^{h-1} \sigma_{\min}(D_k(x)) \sigma_{\min}(W_k). \sigma_{\min}(J_f(x)) \geq \sigma_{\min}(W_h) \prod_{k=1}^{h-1} \sigma_{\min}(D_k(x)) \sigma_{\min}(W_k).$$

And since a diagonal matrix has:

$$\sigma_{\min}(D_k(x)) = \min_i |2u_i(k)(x)|. \sigma_{\min}(D_k(x)) = \min_i |2u_i(k)(x)|.$$


---

## 4. Then inside the trajectory integral we use:

Upper bound:

$$\| J_f(x(t))x'(t) \|_2 \leq \| J_f(x(t)) \|_2 \| x'(t) \|_2. \| J_f(x(t))x'(t) \|_2 \leq \| J_f(x(t)) \|_2 \| x'(t) \|_2.$$

Lower bound:

$$\| J_f(x(t))x'(t) \|_2 \geq \sigma_{\min}(J_f(x(t))) \| x'(t) \|_2. \| J_f(x(t))x'(t) \|_2 \geq \sigma_{\min}(J_f(x(t))) \| x'(t) \|_2.$$

So the two integrals behave oppositely:

$$I(f(x)) = \int_0^1 \| J_f(x(t))x'(t) \| dt. I(f(x)) = \int_0^1 \| J_f(x(t))x'(t) \| dt$$

Upper bound:

$$\leq (\sup_t \| J_f(x(t)) \|_2) I(x) \leq (\sup_t \| J_f(x(t)) \|_2) I(x), I(x) \leq (\sup_t \| J_f(x(t)) \|_2) I(x)$$

Lower bound:

$$\geq (\inf_t \sigma_{\min}(J_f(x(t)))) I(x) \geq (\inf_t \sigma_{\min}(J_f(x(t)))) I(x), I(x) \geq (\inf_t \sigma_{\min}(J_f(x(t)))) I(x)$$

The “directions” of the inequalities come from the fact that one bounds growth and the other bounds shrinkage.

---

## 5. Intuition: Stretching vs squeezing

Think of the Jacobian as a linear transformation that stretches and squeezes space differently in different directions:

### Upper bound (largest singular value)

“How much can this transform stretch the curve in the *worst* direction?”

### Lower bound (smallest singular value)

“How much must this transform stretch the curve in the *least* stretched direction?”

Thus:

- $\sigma_{\max} / \sigma_{\min}$  multiplies into an **upper bound**.
- $\sigma_{\min} / \sigma_{\max}$  multiplies into a **lower bound**.

The inequalities must flip because:

- $\sigma_{\max} / \sigma_{\max}$  represents the **largest** possible magnification factor.
  - $\sigma_{\min} / \sigma_{\min}$  represents the **smallest** possible magnification factor.
- 



## 6. Key summary line

The opposite inequalities arise because:

**Largest singular values control the maximum possible expansion**, so they give “ $\leq$ ” bounds.

**Smallest singular values control the minimum possible expansion (or contraction)**, so they give “ $\geq$ ” bounds.

This is how we get both sides of the trajectory length inequality.

---

If you want, I can write this in LaTeX exactly in the style of your paper’s bound, with the spectral norms and the min/max of the hidden units placed correctly.

You said:

what paper described this notion of expressivity? the topology of  $(h \circ f)^{-1}(e_i) = \{x \in \mathbb{R}^n : h \circ f(x) = e_i\}$  after composing the network with  $h(x) = e_i$  and taking  $x_i = \max_j \{x_j\}$  (Linet et al.,

ChatGPT said:

Good question. There **are** papers that study the *topological* expressivity / capacity of neural networks along exactly those lines; some seem quite related to the idea you describe (preimages, decision-boundaries, topology of level-/preimage-sets, etc.). Below are the main ones, and how well they match your description — with a note on why your “ $(h \circ f)^{-1}(e_i)$  after a max / argmax or pooling-like map” idea may or may not match exactly what they do.

---



## Relevant papers on topological expressivity

### Topology of deep neural networks (by Naitzat, Zhitnikov & Lek-Heng Lim, 2020)

- This paper studies how a data manifold  $M = M_a \cup M_b \subset \mathbb{R}^d = M_a \cup M_b \subset \mathbb{R}^d$  — representing two classes in a binary classification problem — changes topology as it is mapped by a trained neural network. [JMLR+1](#)
- Concretely: They examine how the Betti numbers (topological invariants, roughly counting connected components, holes, loops, etc.) evolve under the network. Their empirical finding: often a complicated topology simplifies drastically, implying the network “untangles” the data topologically. [JMLR](#)
- This is very close in spirit to your notion of looking at preimage sets under a composed map, though they focus on classification and on how data topology evolves, rather than on preimages of “max + indicator” functions per se.

### Topological Expressivity of ReLU Neural Networks (by Ekin Ergen & Moritz Grillo, 2024)

- More recent, this work gives *formal lower and upper bounds* (in terms of architectural parameters like depth/width) on how much a ReLU network can simplify (or change) the topology of data in a binary classification setting. They measure topological changes via Betti numbers. [Proceedings of Machine Learning Research+1](#)

- Their results show that deep networks can perform exponentially greater topological simplification than shallow ones. [Proceedings of Machine Learning Research](#)
- This is perhaps the closest formal work to a theoretical “capacity” or “expressivity” notion defined via topological transformations, which may align with your interest in  $(h \circ f)^{-1}(e_i)(h \circ f)^{-1}(e_i)(h \circ f)^{-1}(e_i)$ -type preimages (though again they do not necessarily use a max/argmax map as “hhh”).

## On The Topological Expressive Power of Neural Networks (by Giovanni Petri & António Leitão, around 2020)

- This earlier work proposes a “topological description of neural network expressive power” in the sense of the space of decision boundaries (for classification) realized by a network architecture. That is, they study how many “different topologies” of decision boundary a network architecture can represent, by sampling across parameter space. [OpenReview](#)
  - Their approach is more abstract: rather than focusing on a particular preimage set for a fixed  $f$ , they consider the space of all decision boundaries the architecture can realize — and measure the “manifold spread” (in some topological sense) of that space. [OpenReview](#)
  - So if your question is about “how expressive is this architecture, in terms of possible topologies of preimage / decision sets,” this is a very relevant paper.
- 

### ⚠ On whether these match exactly “ $(h \circ f)^{-1}(e_i)(h \circ f)^{-1}(e_i)(h \circ f)^{-1}(e_i)$ after max/arg-max / pooling”

- None of the papers above seems to use exactly the same formalism you described (i.e.,  $h(x)=e_i h(x)=e_i$  where  $x_i=\max_j x_j$ ,  $x_i=\max_j x_j$ ). Instead they usually consider classification networks where the final output is a scalar or low-dimensional — and they study the preimage of class-labels / decision-regions (e.g.,  $f^{-1}((-\infty, 0])^f \cap ((-\infty, 0])$  vs  $f^{-1}((0, \infty))^f \cap ((0, \infty))$ ), not necessarily a “max + indicator” map  $h$ .
- The idea of using a “max” or “argmax” layer is more common in architectures with a “max/pooling / maxout” activation, or in soft-max + argmax — but most formal topology-expressivity analyses stick to simpler output codings. For example, the “Topological Expressivity of ReLU Neural Networks” paper is about binary classification

(so effectively a threshold on one real output). [Proceedings of Machine Learning Research+1](#)

- That said: the conceptual motivation behind those works — quantifying how “topologically complicated” a network’s preimage/decision-region capabilities are — is extremely close to what you want. You might try to generalize their techniques, modifying them to the non-standard “max+one-hot” output case (i.e., your  $h \circ f h \circ f$ ).
- 

## My take: the best-fitting prior work & what's novel in your idea

- The 2024 Ergen & Grillo “Topological Expressivity of ReLU Neural Networks” seems like the **state of the art** for formal topological expressivity bounds.
  - The Petri & Leitão “On The Topological Expressive Power of Neural Networks” offers a complementary, architectural-level viewpoint: how expressive a **class** of networks is in terms of possible decision-boundary topologies, rather than for a fixed network instance.
  - Your idea — analyzing the topology of  $(h \circ f)^{-1}(e_i)(h \circ f)^{-1}(e_i)(h \circ f)^{-1}(e_i)$  where  $h$  picks out a “max-class index” — possibly extends those frameworks in an interesting new direction. To my knowledge, **no** paper uses exactly that composition with a “max + one-hot”  $h$  and analyzes the topology of the resulting preimages in a general expressivity bound. That could be a novel angle.
- 

## Recommendation for you

Since your setting (quadratic activations, potentially “max-or-argmax” output) is somewhat nonstandard vs the usual ReLU / binary-classification in the literature, I suggest:

1. Start from the method of **Ergen & Grillo (2024)**: read how they bound topological simplification via Betti numbers. Try to adapt that method to your activation / output mapping.
2. Use the framework of **Petri & Leitão ( $\approx 2020$ )** to think about the space of decision-boundary topologies your architecture can realize — that might give a lower bound on how “diverse” preimage topologies are.

3. Be ready for new technical obstacles: e.g. diagonal Jacobians, non-piecewise-linear activations (quadratic), the fact that “max + one-hot” is not smooth — you may need a variant of their arguments or even new topological tools (stratified sets, singularities, etc.).
- 

If you like — I can **scan the literature (last 5–10 years)** and try to find **all papers** that use *exactly* the kind of preimage-topology / “max-class index + preimage set” expressivity definition you suggest.

Do you want me to build that list for you now?

**You said:**

what are canonical ml datasets of different modes like text and images etc

**ChatGPT said:**

Here's a clean, **mode-by-mode list of canonical machine-learning datasets**, organized the way research papers usually categorize them. I'll include the most historically important, the most widely used benchmarks today, and what each is typically used for.

---



## IMAGE DATASETS

### Classification

- **MNIST** — handwritten digits ( $28 \times 28$  grayscale)
- **Fashion-MNIST** — clothing images (replacement for MNIST)
- **CIFAR-10 / CIFAR-100** — small natural images ( $32 \times 32$  color)
- **ImageNet (ILSVRC)** — 1.2M natural images, 1000 classes (core vision benchmark)
- **TinyImageNet** —  $64 \times 64$  subset of ImageNet (student projects)

### Medical Imaging

- **CheXpert / NIH ChestX-ray14** — chest radiographs

- **ISIC** — dermoscopy images for melanoma detection
- **BraTS** — brain tumor MRI segmentation

## Segmentation / Detection

- **COCO** — object detection, segmentation, captioning
  - **PASCAL VOC** — earlier detection and segmentation benchmark
  - **Cityscapes** — street-scene segmentation (autonomous driving)
  - **ADE20K** — scene parsing
- 



# TEXT / NLP DATASETS

## General NLP tasks

- **Penn Treebank (PTB)** — classic LM benchmark
- **IMDb** — sentiment classification
- **AG News** — news topic classification
- **SST-2 / SST-5** — fine-grained sentiment analysis

## Large-scale language modeling / pretraining

- **WikiText-103** — large Wikipedia LM corpus
- **The Pile** — 825GB curated text dataset (modern foundation-model pretraining)
- **Common Crawl** — web-scale corpus behind GPT/BERT-style pretraining

## Question answering

- **SQuAD v1/v2** — span-extraction QA
- **Natural Questions (NQ)** — Google search QA
- **TriviaQA, HotpotQA** — multi-hop reasoning

## Machine Translation

- **WMT (various years)** — multilingual parallel corpora
  - **IWSLT** — speech and text translation
- 



# AUDIO / SPEECH DATASETS

## Speech Recognition

- **LibriSpeech** — 1000 hours of read English speech
- **TED-LIUM** — TED talk speech transcripts

## Speaker Identification

- **VoxCeleb 1/2**

## Audio Classification

- **ESC-50** — environmental sounds
  - **UrbanSound8K** — real-world environmental audio clips
- 



# VIDEO DATASETS

## Action Recognition

- **Kinetics-400 / Kinetics-700** — large human action video dataset
- **UCF-101** — smaller but widely used
- **HMDB-51** — classic benchmark

## Video object tracking

- **OTB, VOT, LaSOT**

## Video captioning

- **MSR-VTT**
  - **YouCook2**
- 

# TABULAR DATASETS

## Classic ML benchmarks (UCI Repository)

- **Iris** — classic classification
- **Wine Quality**
- **Adult Income**
- **Breast Cancer Wisconsin**
- **Bank Marketing**

## Large-scale tabular datasets

- **KDD Cup 1998/1999**

- **Higgs Boson** dataset
  - **Rossmann Store Sales** (Kaggle)
- 



## MULTIMODAL DATASETS

### Image + Text

- **CLIP's datasets (LAION-400M/5B)** — web image–text pairs
- **COCO Captions** — images with captions
- **Flickr30k** — captioned images
- **Visual Question Answering (VQA)** — pairs of images + natural language questions

### Video + Text

- **HowTo100M** — YouTube instructional videos + transcripts
  - **TVQA** — QA based on TV show clips
- 



## REINFORCEMENT LEARNING DATASETS

### Control / Simulation

- **OpenAI Gym classic control** (CartPole, Acrobot, etc.)
- **MuJoCo** locomotion tasks (HalfCheetah, Ant, Hopper, Walker2d)
- **D4RL** — offline RL benchmark (Maze2D, Adroit, AntMaze)

## Games

- **Atari 2600 suite** (ALE benchmark)
  - **DeepMind Control Suite**
- 



## GRAPHS / NETWORK DATASETS

### Node/graph classification

- **Cora, Citeseer, PubMed** — citation network datasets
- **OGB (Open Graph Benchmark)** — realistic graph datasets (molecules, proteins, citation networks, social networks)

### Molecules

- **QM9** — small molecule quantum properties
  - **ZINC** — graph generation benchmark
- 



## TIME SERIES DATASETS

- **UCR Time Series Archive** — large collection of univariate TS
  - **ECG5000** — ECG classification
  - **Electricity** — forecasting
  - **Traffic** — forecasting
-



# GENERATIVE MODELING DATASETS

- **MNIST, Fashion-MNIST, CIFAR-10** — common for GANs
  - **CelebA** — facial images
  - **LSUN** — large-scale indoor scenes
  - **FFHQ** — high-quality human faces (StyleGAN benchmark)
- 

If you want, I can produce a **short curated list**, a **table**, or a **mode-by-mode “if you’re doing X research, use Y dataset” guide**.