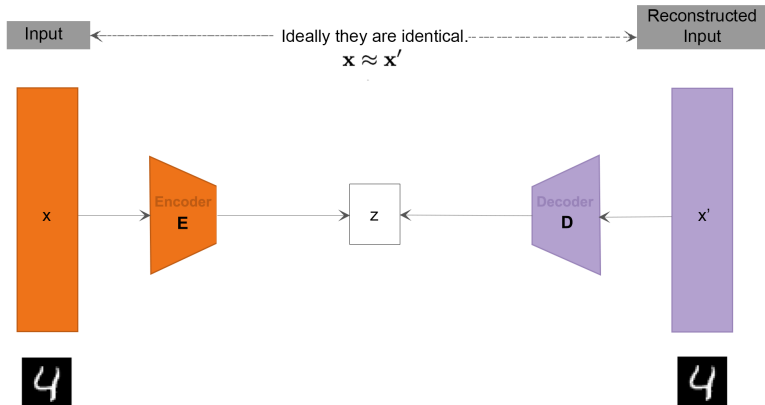


DEEP AUTOENCODERS: FROM UNDERSTANDING TO GENERALIZATION GUARANTEES

Romain Cosentino¹ Randall Balestrieri¹ Richard Baraniuk
¹ Behnaam Aazhang¹

¹Rice University

AUTOENCODERS



AUTOENCODERS: PIECEWISE AFFINE FORMALISM

$$D \circ E(\mathbf{x}) = \sum_{\omega \in \Omega} \mathbf{1}_{\{\mathbf{x} \in \omega\}} (A_{\omega}^D A_{\omega}^E \mathbf{x} + A_{\omega}^D B_{\omega}^E + B_{\omega}^D),$$

- $\mathbf{x} \in \omega \subset \mathbb{R}^d$
- Ω is a partition of the space
- $A_{\omega}^D \in \mathbb{R}^{d \times h}$, $A_{\omega}^E \in \mathbb{R}^{h \times d}$, $B_{\omega}^E \in \mathbb{R}^h$ and $B_{\omega}^D \in \mathbb{R}^d$ with d being the dimension of the input data and h the bottleneck dimension.

$$A_{\omega}^E = W^L Q_{\omega}^{L-1} W^{L-1} \dots Q_{\omega}^1 W^1 \quad \text{and} \quad B_{\omega}^E = \mathbf{b}^L + \sum_{i=1}^{L-1} W^L Q_{\omega}^{L-1} W^{L-1} \dots Q_{\omega}^i \mathbf{b}^i.$$

- $W^{\ell} \in \mathbb{R}^{d_{\ell} \times d_{\ell-1}}$, $\mathbf{b}^{\ell} \in \mathbb{R}^{d_{\ell}}$ the affine parameters of each layer,
- Q^{ℓ} the diagonal matrices encoding the region induced states of the nonlinearities, $(0, 1)$ for ReLU, $(-1, 1)$ for absolute value

2 LAYERS ReLU NETWORK : PIECEWISE AFFINE PARTITIONS

$$f(\mathbf{x}) = W^2 \text{ReLU}(W^1 \mathbf{x}), \text{ where } W^1 = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}, W^2 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}.$$

$$\begin{aligned} \bullet \mathbf{x}_1 &= \begin{pmatrix} 1 \\ 2 \end{pmatrix}, & f(\mathbf{x}_1) &= \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \text{ReLU} \left(\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} \right) \\ & & &= \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}}_{Q_{\omega_1}} \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} \end{aligned}$$

$$\rightarrow \mathbf{x}_1 \in \omega_1$$

$$\begin{aligned} \bullet \mathbf{x}_2 &= \begin{pmatrix} 1 \\ -2 \end{pmatrix}, & f(\mathbf{x}_2) &= \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \text{ReLU} \left(\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} 1 \\ -2 \end{pmatrix} \right) \\ & & &= \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}}_{Q_{\omega_2}} \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} 1 \\ -2 \end{pmatrix} \end{aligned}$$

$$\rightarrow \mathbf{x}_2 \in \omega_2$$

$$\bullet \mathbf{x}_3 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}$$

$$\rightarrow \mathbf{x}_3 \in \omega_1$$

AUTOENCODER INPUT SPACE PARTITIONING

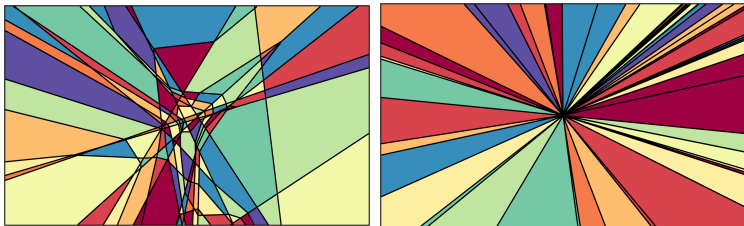


FIGURE 1: 2-dimensional visualizations of the input space partitioning - To reconstruct its input, an AE achieves an affine map for each region - (*Left*) with bias (*Right*) zero bias.

- Each region (described by a specific color) has a particular: $A_{\omega}^D \in \mathbb{R}^{d \times h}$, $A_{\omega}^E \in \mathbb{R}^{h \times d}$, $B_{\omega}^E \in \mathbb{R}^h$ and $B_{\omega}^D \in \mathbb{R}^d$.
- The "code" of each region, $\omega \in \Omega$, is given by the Q_{ω}^{ℓ} .

NUMBER OF REGIONS VS NUMBER OF DATA

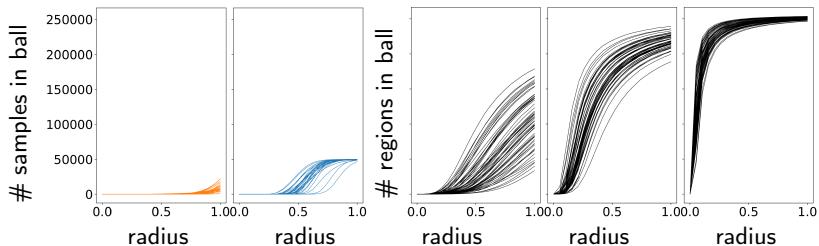


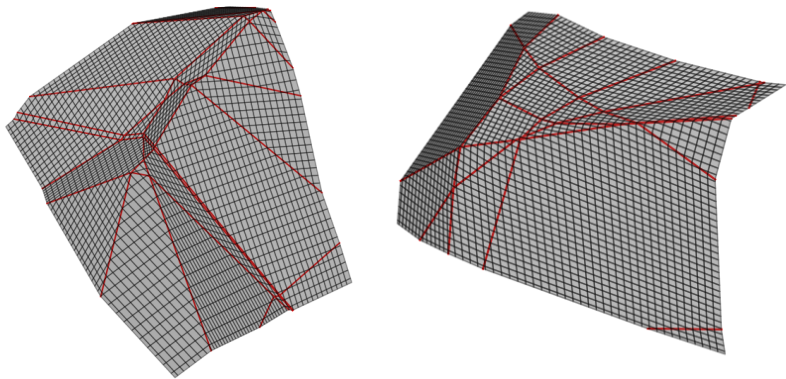
FIGURE 2: *Left:* data - Cifar10 and MNIST. *Right:* regions - small MLP, large MLP, ConvNet.

- Regardless of the dataset and neural network architecture: **the number of regions for any given ball is much larger than the number of data**

OBJECTIVE

Understand how one can control these regions to equip *autoencoders* with **generalisation guarantees**

THE DECODER'S SURFACE



Decoder Continuous Piecewise Affine Surface

Per region Jacobian of the decoder

$$\forall \omega \in \Omega^D, J_\omega[\mathbf{D}] = A_\omega^D,$$

where the columns of A_ω^D form the basis of the tangent space induced by \mathbf{D} .

ASSUMPTION: DATA LIE ON THE ORBIT OF A GROUP



Orbit of digit "7" w.r.t Rotation Group

Which group should we consider?

A LIE GROUP ASSUMPTION

Lie Group: Is a group that is a differentiable manifold.

$$\text{Rotation Group: } SO(2) = \left\{ \begin{pmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{pmatrix} \mid \theta \in \mathbb{R}/2\pi\mathbb{Z} \right\}.$$

Exponential Map: Any matrix Lie group can be defined via an exponential map.

$$\text{Rotation Group: } SO(2) = \left\{ \exp(\theta G) \mid G = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}, \theta \in \mathbb{R}/2\pi\mathbb{Z} \right\}.$$

Orbit w.r.t Lie Group: The data x are modeled by

$$x(\theta) = \exp(\theta G)x(0),$$

the orbit of $x(0)$ with respect to the group induced by $\exp(\theta G)$.

EQUIVARIANCE LIE GROUP REGULARIZATIONS

$$\mathcal{L}_{\text{LieReg}} = \underbrace{\sum_{i=1}^n \|\mathbf{D}(\mathbf{E}(x_i)) - x_i\|}_{\text{Reconstruction Error}} + \underbrace{\min_G \sum_{\omega \in \Omega^D} \min_{\theta} \|\exp(\theta G) A_{\omega_0}^D - A_{\omega}^D\|}_{\text{Lie Group Assumption}}$$

Learning with an exponential map

$$\exp(\theta G)$$

1. Non-Convex
2. Tedious Computation of the Gradient

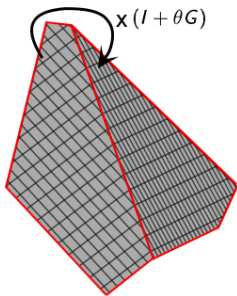
Local approximation

$$\exp(\theta G) \approx_{\theta \sim 0} I + \theta G$$

1. Only for small transformations
2. Need to know the neighbors of each sample

LOCALLY CONSTRAINING THE CURVATURE

$$\mathcal{L}_{\text{LieReg}} \approx \sum_{i=1}^n \|\mathbf{D}(\mathbf{E}(x_i)) - x_i\| + \min_G \sum_{\omega \in \Omega^D} \sum_{\omega' \in \mathcal{N}(\omega)} \min_{\theta} \left\| A_{\omega}^D - (I + \theta G) A_{\omega'}^D \right\|,$$



GENERALISATION GUARANTEE

1. **If** the Decoder approximates the tangent space of the data at a position.
2. **If** the Lie group regularization is 0.

Then

The approximation of the data manifold is upper-bounded by the sum of the radius of each region.

THEOREM

If on a region $\omega' \in \Omega^D$ the matrix $A_{\omega'}^D$ forms a basis of the manifold tangent space on this region, and the Lie group regularization is 0 then for all region $\omega \in \Omega^D$ the basis vectors of A_{ω}^D are the basis vector of the tangent of the data manifold with

$$d(\cup_{\omega \in \Omega^D} \mathcal{T}_{AE}(\omega), \mathcal{X}) \leq \sum_{\omega \in \Omega^D} \text{Rad}(\omega),$$

where $\mathcal{T}_{AE}(\omega)$ the tangent space of the AE for the region ω , \mathcal{X} denotes the data manifold, d defines the 2-norm distance, and $\text{Rad}(\omega_i)$ the radius of the region ω_i .

RESULTS

TABLE 1: Comparison of the testing reconstruction errors ($\times 10^{-2} \pm \text{std} \times 10^{-2}$)

<i>Dataset</i> <i>Model</i>	AE	Den. AE	H.O.C. AE	Lie Group
CIFAR10	5.6 \pm 0.05	5.0 \pm 0.05	-	4.9 \pm 0.07
MNIST	12.01 \pm 0.003	12.01 \pm 0.004	12.01 \pm 0.004	6.3 \pm 0.1
CBF	62.38 \pm 0.74	52.66 \pm 0.76	51.09 \pm 0.54	43.99 \pm 1.2
Yoga	33.76 \pm 0.81	33.29 \pm 0.72	32.08 \pm 0.42	20.28 \pm 1.1
Trace	13.95 \pm 0.45	11.28 \pm 0.57	12.57 \pm 0.21	10.91 \pm 0.45
Wine	63.06 \pm 0.02	59.34 \pm 0.02	49.94 \pm 0.02	19.01 \pm 0.02
ShapesAll	67.98 \pm 3.0	58.67 \pm 1.4	61.42 \pm 5.5	52.97 \pm 1.9
FiftyWords	64.91 \pm 1.7	60.91 \pm 1.0	60.92 \pm 0.7	57.89 \pm 1.0
WordSynonyms	70.95 \pm 1.5	66.02 \pm 0.8	66.52 \pm 0.5	62.22 \pm 1.1
InsectSounds	51.86 \pm 0.6	40.24 \pm 0.8	41.93 \pm 0.6	38.11 \pm 0.9
ECG5000	21.92 \pm 0.75	20.31 \pm 0.39	20.31 \pm 0.36	18.06 \pm 0.9
Earthquakes	56.23 \pm 4.1	54.62 \pm 4.1	51.79 \pm 1.0	50.20 \pm 0.5
Haptics	37.25 \pm 0.2	36.02 \pm 1.8	27.21 \pm 0.5	16.94 \pm 3.4
FaceFour	49.82 \pm 1.0	48.51 \pm 0.8	48.52 \pm 0.7	46.00 \pm 0.6
Synthetic	95.61 \pm 1.3	89.37 \pm 1.0	88.47 \pm 0.9	55.87 \pm 0.8

CONCLUSION & DIRECTIONS

- We propose a way to learn an equivariant AE.
 - The underlying group is learned via the Lie group generator G .
 - Under the Lie group assumption on the data, we obtain generalization guarantees.
 - Propose a way to develop constraints on the approximated manifold that are assumption driven.
-

- Learning Lie group is non trivial.
- Generalizing to "pancakes" and multiple orbits.
- Provide efficient and principled ways to sample neighboring regions.