

Neural Network Approximation based on Hausdorff distance of Tropical Zonotopes

Panagiotis Misiakos, George Smyrnis, George Retsinas, Petros Maragos



National Technical University of Athens
School of Electrical and Computer Engineering

International Conference on Learning Representations (ICLR) 2022

Contributions

- ✓ **Novel** bound on neural network approximation.
- ✓ 2 **new** algorithms for neural network compression.

Tropical Algebra

Tropical Geometry

Tropical Algebra

✓ **Tropical Semiring** $\mathbb{R}_{\max} = \mathbb{R} \cup \{-\infty\}$

$$a \vee b = \max(a, b)$$

$$a + b = a + b$$

- Replaces classical operations of addition and multiplication with max and +, respectively.

Tropical Geometry

Tropical Algebra

✓ **Tropical Semiring** $\mathbb{R}_{\max} = \mathbb{R} \cup \{-\infty\}$

$$a \vee b = \max(a, b)$$

$$a + b = a + b$$

- Replaces classical operations of addition and multiplication with \max and $+$, respectively.

✓ **Tropical Polynomials**

$$f(\mathbf{x}) = \max_{i \in [n]} \{\mathbf{a}_i^T \mathbf{x} + b_i\}$$

- Expressive for ReLU networks.

Tropical Geometry

Tropical Algebra

✓ **Tropical Semiring** $\mathbb{R}_{\max} = \mathbb{R} \cup \{-\infty\}$

$$a \vee b = \max(a, b)$$

$$a + b = a + b$$

- Replaces classical operations of addition and multiplication with \max and $+$, respectively.

✓ **Tropical Polynomials**

$$f(\mathbf{x}) = \max_{i \in [n]} \{\mathbf{a}_i^T \mathbf{x} + b_i\}$$

- Expressive for ReLU networks.

Tropical Geometry

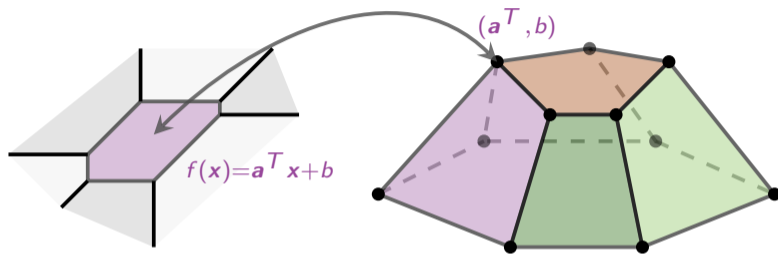
✓ **Newton Polytopes**

$$\text{Newt}(f) = \text{conv} \{\mathbf{a}_i : i \in [n]\}$$

$$\text{ENewt}(f) = \text{conv} \{(\mathbf{a}_i, b_i) : i \in [n]\}$$

- They provide geometric interpretation for tropical polynomials.

Linear Regions and the Newton Polytope



- ✓ 1 – 1 mapping: between linear regions and vertices. [1]
- ✓ The upper envelope determines the tropical polynomial and vice versa

$$f, g \in \mathbb{R}_{\max}[\mathbf{x}] : f = g \Leftrightarrow UF(\text{ENewt}(f)) = UF(\text{ENewt}(g))$$

[1] Charisopoulos, V., Maragos, P. A tropical approach to neural networks with piecewise linear activations. *arXiv preprint arXiv:1805.08749*, 2018

Idea: What if we relax the previous equality?

Question: Would $\text{ENewt}(f) \approx \text{ENewt}(g)$ imply $f \approx g$?

Idea: What if we relax the previous equality?

Question: Would $\text{ENewt}(f) \approx \text{ENewt}(g)$ imply $f \approx g$?

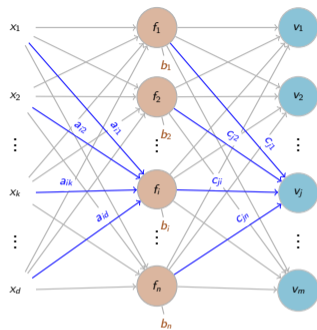
Proposition

Let $p, \tilde{p} \in \mathbb{R}_{\max}[\mathbf{x}]$ be two tropical polynomials and let $P = \text{ENewt}(p)$, $\tilde{P} = \text{ENewt}(\tilde{p})$. Then,

$$\max_{\mathbf{x} \in \mathcal{B}} |p(\mathbf{x}) - \tilde{p}(\mathbf{x})| \leq \rho \cdot \mathcal{H}(P, \tilde{P})$$

where $\mathcal{B} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| \leq r\}$ is the hypersphere of radius r , and $\rho = \sqrt{r^2 + 1}$.

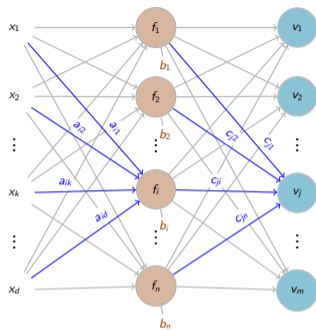
Tropical Geometry of Neural Networks



ReLU neural network with 1 hidden layer

[2] L. Zhang, G. Naitzat, L.-H. Lim. "Tropical Geometry of Deep Neural Networks." in *International Conference on Machine Learning*, pages 5824–5832. 2018.

[3] P. Maragos, V. Charisopoulos and E. Theodosis, "Tropical Geometry and Machine Learning," in *Proceedings of the IEEE*, vol. 109, no. 5, pp. 728-755, May 2021, doi: 10.1109/JPROC.2021.3065238.



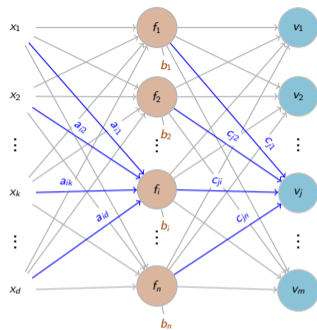
ReLU neural network with 1 hidden layer

✓ i -th hidden layer node.

$$f_i(\mathbf{x}) = \max(\mathbf{a}_i^T \mathbf{x} + b_i, 0)$$

[2] L. Zhang, G. Naitzat, L.-H. Lim. "Tropical Geometry of Deep Neural Networks." in *International Conference on Machine Learning*, pages 5824–5832. 2018.

[3] P. Maragos, V. Charisopoulos and E. Theodosis, "Tropical Geometry and Machine Learning," in *Proceedings of the IEEE*, vol. 109, no. 5, pp. 728-755, May 2021, doi: 10.1109/JPROC.2021.3065238.



ReLU neural network with 1 hidden layer

✓ i -th hidden layer node.

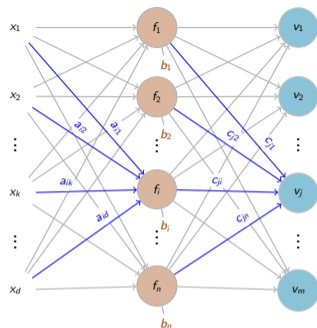
$$f_i(\mathbf{x}) = \max(\mathbf{a}_i^T \mathbf{x} + b_i, 0)$$

✓ j -th output node.

$$v_j(\mathbf{x}) = p_j(\mathbf{x}) - q_j(\mathbf{x})$$

[2] L. Zhang, G. Naitzat, L.-H. Lim. "Tropical Geometry of Deep Neural Networks." in *International Conference on Machine Learning*, pages 5824–5832. 2018.

[3] P. Maragos, V. Charisopoulos and E. Theodosis, "Tropical Geometry and Machine Learning," in *Proceedings of the IEEE*, vol. 109, no. 5, pp. 728-755, May 2021, doi: 10.1109/JPROC.2021.3065238.



ReLU neural network with 1 hidden layer

- ✓ i -th hidden layer node.

$$f_i(\mathbf{x}) = \max(\mathbf{a}_i^T \mathbf{x} + b_i, 0)$$

- ✓ j -th output node.

$$v_j(\mathbf{x}) = p_j(\mathbf{x}) - q_j(\mathbf{x})$$

Tropical Geometry

- ✓ $\text{ENewt}(f_i)$ is linear segment with endpoints $\mathbf{0}$ and (\mathbf{a}_i^T, b_i) .
- ✓ $P_j = \text{ENewt}(p_j)$, $Q_j = \text{ENewt}(q_j)$ are Minkowski sums of segments \Leftrightarrow **zonotopes** [2,3].
- ✓ (\mathbf{a}_i^T, b_i) are called **generators**.

[2] L. Zhang, G. Naitzat, L.-H. Lim. "Tropical Geometry of Deep Neural Networks." in *International Conference on Machine Learning*, pages 5824–5832. 2018.

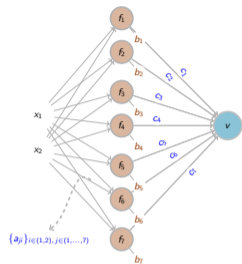
[3] P. Maragos, V. Charisopoulos and E. Theodosis, "Tropical Geometry and Machine Learning," in *Proceedings of the IEEE*, vol. 109, no. 5, pp. 728-755, May 2021, doi: 10.1109/JPROC.2021.3065238.

Theorem

Let $v, \tilde{v} \in \mathbb{R}_{\max}[\mathbf{x}]$ be two neural networks with 1 hidden layer and \tilde{P}_j, \tilde{Q}_j denote the positive and negative zonotopes of \tilde{v} . The following bound applies.

$$\max_{\mathbf{x} \in \mathcal{B}} \|v(\mathbf{x}) - \tilde{v}(\mathbf{x})\|_1 \leq \rho \cdot \left(\sum_{j=1}^m \mathcal{H}(P_j, \tilde{P}_j) + \mathcal{H}(Q_j, \tilde{Q}_j) \right)$$

- ✓ Geometrical approximation problem.
- ✓ **Goal:** approximate the zonotopes.

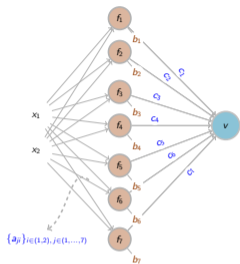


(a) Original network

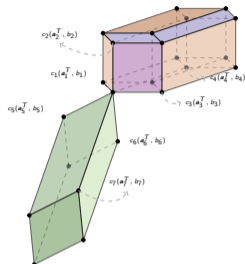
✓ Applies only to networks with one output neuron.

Zonotope K-means

Compression Algorithms I. Zonotope K-means



(a) Original network



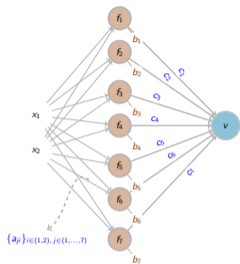
(b) Original zonotopes

✓ Applies only to networks with one output neuron.

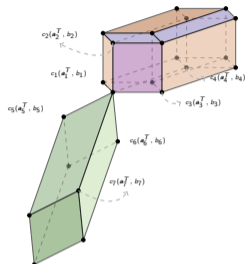
Zonotope K-means

1. Split zonotope generators into positive and negative.

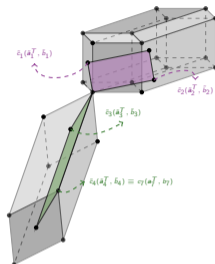
Compression Algorithms I. Zonotope K-means



(a) Original network



(b) Original zonotopes



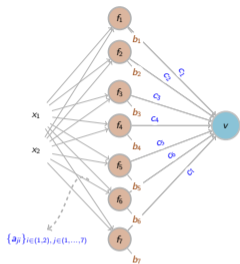
(c) Resulting zonotopes.

✓ Applies only to networks with one output neuron.

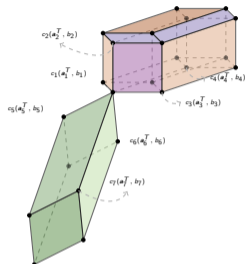
Zonotope K-means

1. Split zonotope generators into positive and negative.
2. Apply K-means to each generating set.

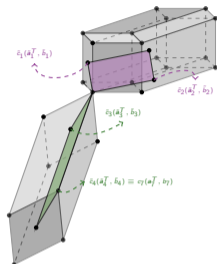
Compression Algorithms I. Zonotope K-means



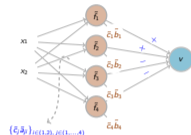
(a) Original network



(b) Original zonotopes



(c) Resulting zonotopes.

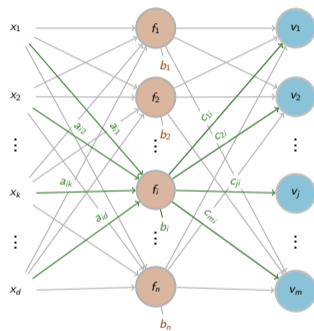


(d) Compressed network.

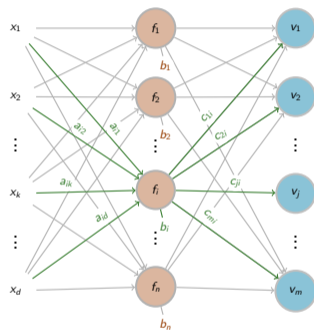
✓ Applies only to networks with one output neuron.

Zonotope K-means

1. Split zonotope generators into positive and negative.
2. Apply K-means to each generating set.
3. Construct final network.

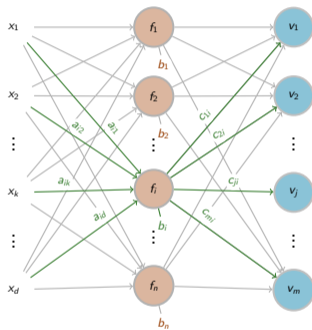


- Zonotope K-means doesn't generalize directly to multiple output nodes.



- Zonotope K-means doesn't generalize directly to multiple output nodes.

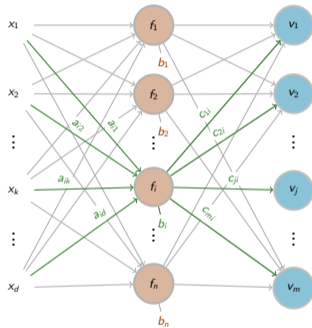
Neural Path K-means



- Zonotope K-means doesn't generalize directly to multiple output nodes.

Neural Path K-means

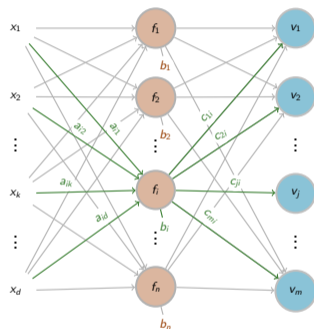
1. For each node form the vector of weights of incident edges.



- Zonotope K-means doesn't generalize directly to multiple output nodes.

Neural Path K-means

1. For each node form the vector of weights of incident edges.
2. Execute K-means to these vectors.



- Zonotope K-means doesn't generalize directly to multiple output nodes.

Neural Path K-means

1. For each node form the vector of weights of incident edges.
2. Execute K-means to these vectors.
3. Construct reduced network.

Zonotope K-means Bound

$$\frac{1}{\rho} \cdot \max_{\mathbf{x} \in \mathcal{B}} |v(\mathbf{x}) - \tilde{v}(\mathbf{x})| \leq K \cdot \delta_{\max} + \left(1 - \frac{1}{N_{\max}}\right) \sum_{i=1}^n |c_i| \left\| \left(\mathbf{a}_i^T, b_i\right) \right\|$$

Neural Path K-means Bound

$$\begin{aligned} \frac{1}{\rho} \cdot \max_{\mathbf{x} \in \mathcal{B}} \|v(\mathbf{x}) - \tilde{v}(\mathbf{x})\|_1 &\leq \sqrt{m} K \delta_{\max}^2 + \sqrt{m} \left(1 - \frac{1}{N_{\max}}\right) \sum_{i=1}^n \|C_{:,i}\| \left\| \left(\mathbf{a}_i^T, b_i\right) \right\| + \\ &\frac{\sqrt{m} \delta_{\max}}{N_{\min}} \sum_{i=1}^n \left(\left\| \left(\mathbf{a}_i^T, b_i\right) \right\| + \|C_{:,i}\| \right) + \sum_{j=1}^m \sum_{i \in \mathcal{N}_j} |c_{ji}| \left\| \left(\mathbf{a}_i^T, b_i\right) \right\| \end{aligned}$$

- ✓ Theoretical bounds depend on the weights magnitude and K-means parameters.
- ✓ Approximation is better when $K \approx n$. Both bounds become 0 when $K = n$.

Experimental Evaluation I: Comparison with tropical techniques.

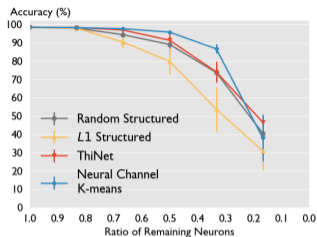
✓ Binary classification tasks.

Percentage of Remaining Neurons	MNIST 3/5			MNIST 4/9		
	Smyrnis et al., 2020	Zonotope K-means	Neural Path K-means	Smyrnis et al., 2020	Zonotope K-means	Neural Path K-means
100% (Original)	99.18 ± 0.27	99.38 ± 0.09	99.38 ± 0.09	99.53 ± 0.09	99.53 ± 0.09	99.53 ± 0.09
5%	99.12 ± 0.37	99.42 ± 0.07	99.25 ± 0.04	98.99 ± 0.09	99.52 ± 0.09	99.48 ± 0.15
1%	99.11 ± 0.36	99.39 ± 0.05	99.32 ± 0.03	99.01 ± 0.09	99.46 ± 0.05	99.35 ± 0.17
0.5%	99.18 ± 0.36	99.41 ± 0.05	99.22 ± 0.11	98.81 ± 0.09	99.35 ± 0.24	98.84 ± 1.18
0.3%	99.18 ± 0.36	99.25 ± 0.37	99.19 ± 0.41	98.81 ± 0.09	98.22 ± 1.38	98.22 ± 1.33

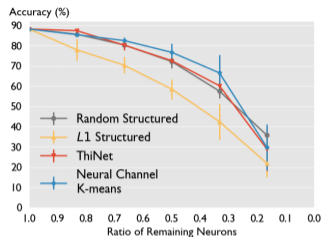
✓ Multiclass classification tasks.

Percentage of Remaining Neurons	MNIST		Fashion-MNIST	
	Smyrnis and Maragos, 2020	Neural Path K-means	Smyrnis and Maragos, 2020	Neural Path K-means
100% (Original)	98.60 ± 0.03	98.61 ± 0.11	88.66 ± 0.54	89.52 ± 0.19
50%	96.39 ± 1.18	98.13 ± 0.28	83.30 ± 2.80	88.22 ± 0.32
25%	95.15 ± 2.36	98.42 ± 0.42	82.22 ± 2.85	86.67 ± 1.12
10%	93.48 ± 2.57	96.89 ± 0.55	80.43 ± 3.27	86.04 ± 0.94
5%	92.93 ± 2.59	96.31 ± 1.29	–	83.68 ± 1.06

Experimental Evaluation II: Comparison with Thinet and baselines.



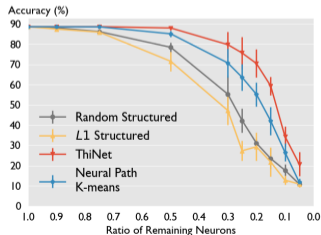
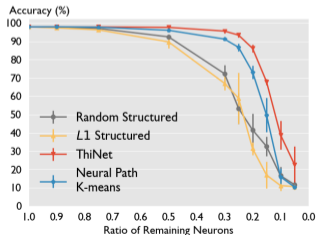
(a) LeNet5, MNIST



(b) LeNet5, F-MNIST

LeNet5

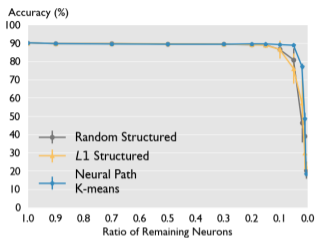
- ✓ 1 hidden layer with 84 neurons.



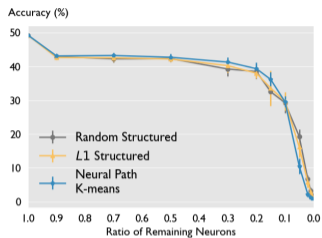
Custom deep network

- ✓ 3 hidden layers.

Experimental Evaluation III: Larger datasets



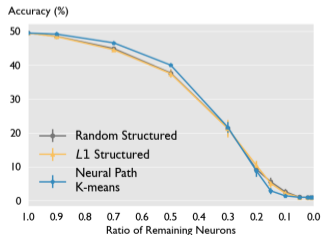
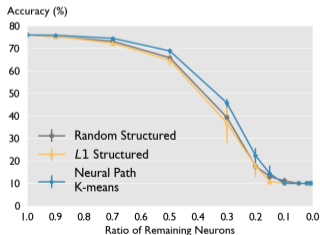
(a) CIFAR-VGG, CIFAR10



(b) CIFAR-VGG, CIFAR100

CIFAR-VGG

- ✓ 1 hidden layer of size 512.



AlexNet

- ✓ 2 hidden layers of size 512.