

# Certifying Differential Invariants of Neural Networks using Abstract Duals

CHANDRA KANTH NAGESH, University of Colorado Boulder, USA

Neural networks are increasingly used in safety-critical domains, necessitating formal guarantees about their properties under perturbations to inputs. Existing robust verification techniques, typified by DeepPoly, primarily focus on the robustness analysis during forward mode learning where certifying output stability is typified by  $f(I) \subseteq [y_L, y_R]$  for an  $L_\infty$ -constrained input set  $I = \{x \mid \|x - x_0\|_\infty \leq \epsilon\}$ . While such methods employ a sophisticated polyhedral abstract domain (combining intervals and affine forms) to generate sound, tight bounds, this entire class of analysis still fails to provide sound guarantees over the behaviour of the network. This oversight creates a critical verification gap for gradient-dependent systems where the Backward Pass Verification Problem which can be described formally as bounding the Jacobian  $J(x) = \nabla_x f(x)$  such that  $J(x) \subseteq [J_L, J_R]$  for all  $x \in I$  is essential.

## 1 Introduction

DeepPoly[1]

## 2 Overview

### 3 [Contribution 1]

### 4 [Contribution 2]

## 5 Evaluation

## 6 Related Work

## 7 Conclusion

## Acknowledgments

TBD

## References

- [1] Gagandeep Singh, Timon Gehr, Markus Püschel, and Martin Vechev. 2019. An abstract domain for certifying neural networks, Vol. 3. 1–30. doi:10.1145/3290354

---

Author's Contact Information: Chandra Kanth Nagesh, ckn@colorado.edu, University of Colorado Boulder, Boulder, Colorado, USA.