

Training Robust Neural Networks

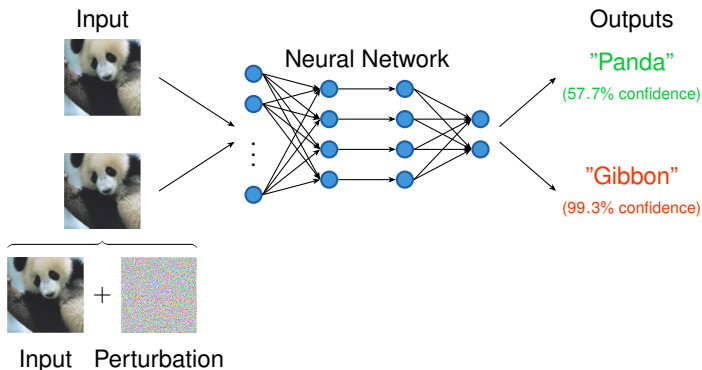
Lukas Koller

Prof. Dr.-Ing. Matthias Althoff
Cyber-Physical Systems Group
Technische Universität München

February 4, 2024

Neural Networks and Adversarial Attacks

Adversarial Attack: Small carefully chosen input perturbation that leads to a misclassification.

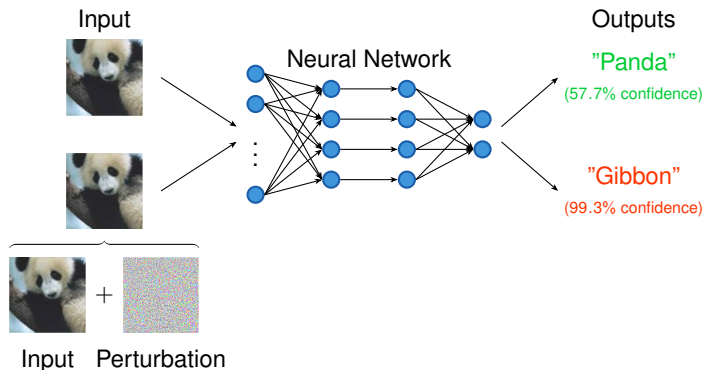


Neural networks are vulnerable to adversarial attacks!¹

¹Goodfellow et al., "Explaining and Harnessing Adversarial Examples"

Neural Networks and Adversarial Attacks

Adversarial Attack: Small carefully chosen input perturbation that leads to a misclassification.



Neural networks are vulnerable to adversarial attacks!¹

¹Goodfellow et al., "Explaining and Harnessing Adversarial Examples"

Training Robust Neural Networks

There are many different methods to train adversarially robust neural networks, e.g. ², ³, ⁴, ⁵.

Your tasks: Review and compare different approaches to train and evaluate robust neural networks.

Interested? Contact me!

Lukas Koller
lukas.koller@tum.de

²Madry et al., “Towards Deep Learning Models Resistant to Adversarial Attacks”

³Gowal et al., “Scalable Verified Training for Provably Robust Image Classification”

⁴Mirman et al., “Differentiable Abstract Interpretation for Provably Robust Neural Networks”

⁵Zhang et al., “Theoretically Principled Trade-off between Robustness and Accuracy”