# Boosting Model Robustness by Leveraging Data Augmentations, Stability Training, and Noise Injections
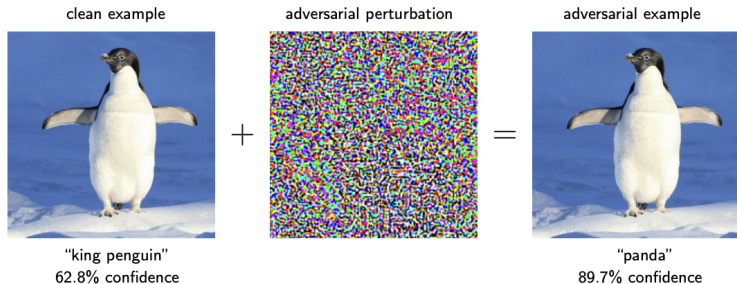
Soon Hoe Lim (Nordita, KTH & Stockholm U)

Joint work with N. Benjamin Erichson (U of Pittsburgh), Francisco Utrera (ICSI & U of Pittsburgh), Winnie Xu (U of Toronto), Ziang Cao (U of Pittsburgh), and Michael Mahoney (ICSI & UC Berkeley)

**Applied Machine Learning Days 2022**

# Model Robustness Matters!

- Deep learning models are typically brittle and sensitive to noisy and adversarial environments
- For many real-world applications, obtaining stable and robust statistical performance is more important than simply achieving SOTA predictive performance
- Here we focus on input stability (robustness) with respect to common data corruptions and domain shifts that naturally occur in many real-world applications



clean example
"king penguin"
62.8% confidence

adversarial perturbation

adversarial example
"panda"
89.7% confidence

▶ Szegedy et al. "Intriguing properties of neural networks." ICLR (2014).
▶ Goodfellow et al. "Explaining and harnessing adversarial examples." ICLR (2015).
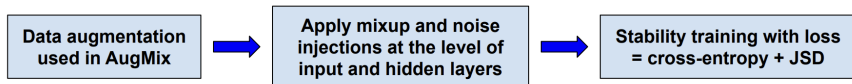
## NoisyMix

Four common methods to improve model robustness to input perturbations are:

- Data augmentations
- Stability training
- Mixup
- Noise injections

**How can we leverage the strength of these methods to further improve both model robustness and test accuracy?**
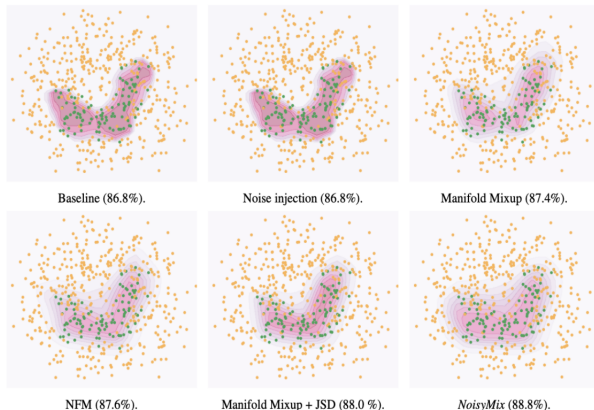
We introduce NoisyMix, a training scheme that judiciously combines all of the above components in a single setup to boost both robustness and accuracy in classification tasks.

Sounds simple, but the devil is in the details to get the method right:

| Data augmentation used in AugMix | → | Apply mixup and noise injections at the level of input and hidden layers | → | Stability training with loss = cross-entropy + JSD |
|---|---|---|---|---|

# Diving Deeper into NoisyMix

The advantage of NoisyMix compared to other schemes is illustrated on a binary classification task on a noisy toy dataset (without augmentation), where it can be seen that NoisyMix is most effective at smoothing the decision boundary and yields the best test accuracy:
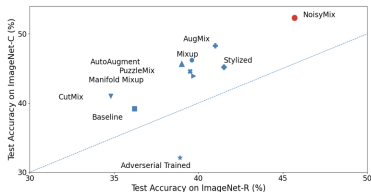


Baseline (86.8%).

Noise injection (86.8%).

Manifold Mixup (87.4%).

NFM (87.6%).

Manifold Mixup + JSD (88.0 %).

*NoisyMix* (88.8%).

Moreover, we provide theory to understand the effects of NoisyMix through the lens of implicit regularization and show that minimizing the NoisyMix loss can lead to a small regularized adversarial loss and a stable model (see paper).

# Empirical Results

We benchmark common corruptions with **RobustBench**, and find that NoisyMix (currently) tops the leaderboards (CIFAR-10-C, CIFAR-100-C, and ImageNet-C) there.



More results are available at:

- Paper: `https://arxiv.org/abs/2202.01263`
- RobustBench: `https://robustbench.github.io/`