

Robustness of deep networks

For image classifiers, achieving robustness to perturbations is a crucial requirement



Lampshade

E.g. adversarial perturbations causing misclassification



Universal adversarial perturbation

A single perturbation causing misclassification with high probability



Universal adversarial perturbations

Seyed-Mohsen Moosavi-Dezfooli*, Alhussein Fawzi+, Omar Fawzi+, Pascal Frossard* *EPFL, +UCLA, +ENS-Lyon

The algorithm

Finding universal perturbation using the set X



Algorithm Computation of universal perturbations.

- 1: **input:** Data points X, classifier \hat{k} , desired ℓ_p norm of the perturbation ξ , desired accuracy on perturbed samples δ .
- 2: **output:** Universal perturbation vector v.
- 3: Initialize $v \leftarrow 0$.
- 4: while $\operatorname{Err}(X_v) \leq 1 \delta$ do
- 5: **for** each datapoint $x_i \in X$ do
- 6: if $\hat{k}(x_i + v) = \hat{k}(x_i)$ then
- 7: Compute the **minimal** perturbation that sends $x_i + v$ to the decision boundary:

$$\Delta v_i \leftarrow \arg\min_{x} ||r||_2 \text{ s.t. } \hat{k}(x_i + v + r) \neq \hat{k}(x_i).$$

8:

Update the perturbation:

 $v \leftarrow \mathcal{P}_{p,\xi}(v + \Delta v_i).$

9: end if
10: end for
11: end while

Misclassification rate

High misclassification rate on unseen test data



What is special in these perturbations?

Comparison of universal perturbation to other types of perturbations

Fooling percentage of the perturbations with the same norm



Random noise achieving 90% fooling rate









Doubly universal perturbations

The ability to generalize well across different neural networks

Cross-model fooling percentage for VGG-16







Perturbed labels

Existence of *dominant* labels that suck most other labels.



High volume classification regions!

Feedbacking does not help!

Augmenting the training data with perturbed images



... because universal perturbations are diverse!



Explaining universal perturbations

The decision boundary in the vicinity of natural images is highly correlated!



Check our new work: "Analysis of universal adversarial perturbations"

github.com/lts4/universal