# Adversarial Label Flips
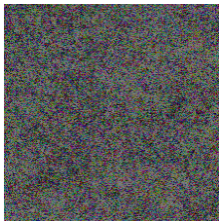
Matthias Dellago & Maximilian Samsinger

A short recap
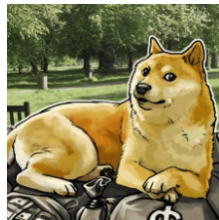
# Adversarial attack



Husky
(42.82% confidence)

$+\ \epsilon$

Noise (PGD-40)
50x amplified

$=$

Handkerchief
(99.999988% confidence)

Source: `ctf.codes`, circa 2021

# What we want to do

**Confusion Matrix**

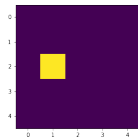|  |  | Categorised as | | |
|---|---|---|---|---|
|  |  | Dog | Cat | Plane |
|  | Dog | 0.0 | ? | ? |
| Adversarial Example of a | Cat | ? | 0.0 | ? |
|  | Plane | ? | ? | 0.0 |

How many modified dogs get classified as cats vs as planes? etc.

Some simple theory

We want similar images that are classified differently.
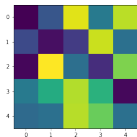But what is "similar"?

# Quantifying Changes

| $L^0$-Norm | $L^1$-Norm | $L^2$-Norm | $L^\infty$-Norm |
|---|---|---|---|
| Number of pixels changed | Sum of all changes | Sum of the *square* of all changes | Maximum of all changes |



| Perturb one pixel maximally | Minimise sum | Minimise sum of squares | Perturb all pixels equally |
|---|---|---|---|

# Two Different Approaches

small displacement and
big misclassification

maximise misclassification

minimise displacement

small displacement

big misclassification

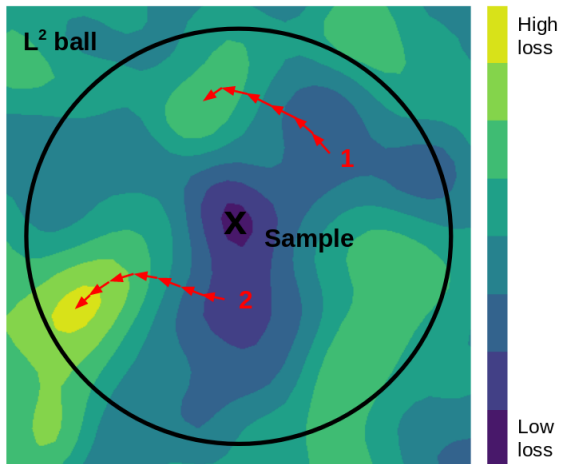# Projected Gradient Decent

1. Pick spot in epsilon ball
2. Iterate gradient decent
3. If leaving ball, project back onto surface.



Towards Deep Learning Models Resistant to Adversarial Attacks, Aleksander Madry et al., arXiv, 2019

# Projected Gradient Decent



Know your enemy, Oscar Knagg, towardsdatascience.com, 2019

universität
innsbruck

# Carlini-Wagner-Attack

Original approach: minimise disance while always staying in "misclassification territory".

Problem: Nonlinearity of constraint makes for bad optimisation properties.

Towards Evaluating the Robustness of Neural Networks, Nicholas Carlini and David Wagner, IEEE, 2017

# Carlini-Wagner-Attack

Solution: Pack constraint into the function that is optimised.

$\rightarrow$ minimise: distance - "how misclassified is x?"*

i.e. minimise distance while maximising misclassification.

*loss function

Towards Evaluating the Robustness of Neural Networks, Nicholas Carlini and David Wagner, IEEE, 2017
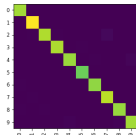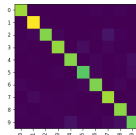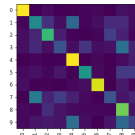
# Code

# Results

# MNIST, $L^\infty$-PGD

# MNIST, $L^2$-Carlini-Wagner-Attack

# CIFAR-10, $L^\infty$-PGD

# CIFAR-10, $L^0$-Brendel-Bethge-Attack

# References I

📄 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus.
Intriguing properties of neural networks.
In *International Conference on Learning Representations (ICLR)*, 2014.

📄 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy.
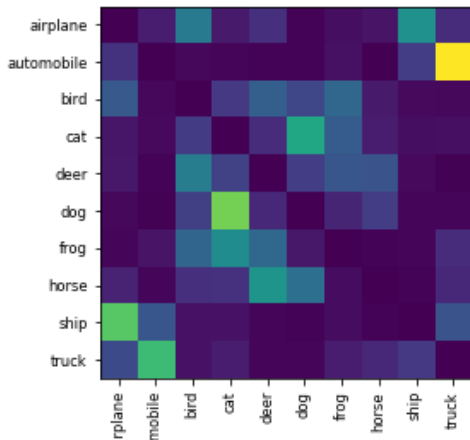Explaining and harnessing adversarial examples.
*arXiv preprint arXiv:1412.6572*, 2014.

📄 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
Towards deep learning models resistant to adversarial attacks.
*arXiv preprint arXiv:1706.06083*, 2017.

# References II

📄 Jonas Rauber, Wieland Brendel, and Matthias Bethge.
Foolbox: A python toolbox to benchmark the robustness of
machine learning models.
*arXiv preprint arXiv:1707.04131*, 2017.

📄 Li Deng.
The mnist database of handwritten digit images for machine
learning research [best of the web].
*IEEE Signal Processing Magazine*, 29(6):141–142, 2012.

📄 Han Xiao, Kashif Rasul, and Roland Vollgraf.
Fashion-mnist: a novel image dataset for benchmarking machine
learning algorithms.
*arXiv preprint arXiv:1708.07747*, 2017.

# References III

Alex Krizhevsky, Geoffrey Hinton, et al.
Learning multiple layers of features from tiny images.
2009.