# Adversarial Label Flips
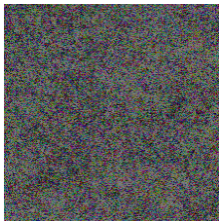
Matthias Dellago & Maximilian Samsinger

A short recap
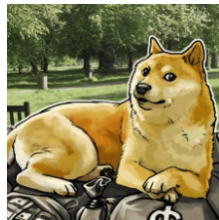
# Fast gradient sign method



Husky
(42.82% confidence)

Noise (PGD-40)
50x amplified

Handkerchief
(99.999988% confidence)

Source: `ctf.codes`, circa 2021

# What we want to do

## Confusion Matrix

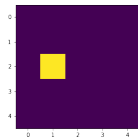|  |  | Categorised as | | |
|---|---|---|---|---|
|  |  | Dog | Cat | Plane |
| Adversarial Example of a | Dog | 0.0 | ? | ? |
|  | Cat | ? | 0.0 | ? |
|  | Plane | ? | ? | 0.0 |

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"?
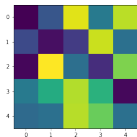
# Quantifying Changes

| $L^0$-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 |
|---|---|---|---|

universität
innsbruck

# Two Different Approaches

small displacement and
big misclassification

maximise misclassification

minimise displacement

small displacement

big misclassification
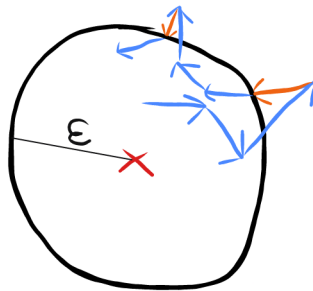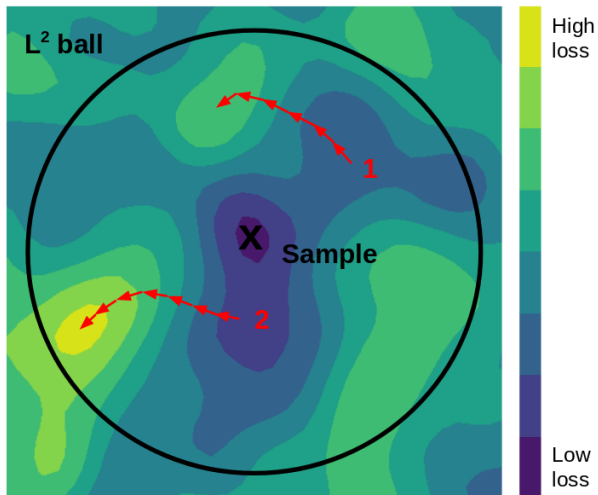
# Projected Gradient Decent

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

# Projected Gradient Decent



L² ball

**X** Sample

1

2

High loss

Low loss

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

# Carlini-Wagner-Attack

Original idea: minimise disance while always staying in
"misclassification territory".
Problem: Nonlinearity of constraint makes for bad optimisation

# Carlini-Wagner-Attack

Solution: Pack constraint into the function that is optimised.
=> minimise: distance + "how misclassified is x?"*
equiv. minimise distance while maximising misclassification.

*loss function

# References I