

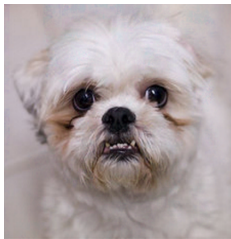
# Adversarial Examples

A new evil has announced its arrival...

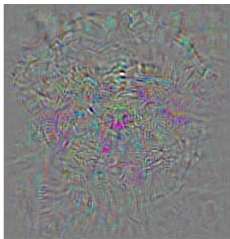
Sewade Ogun

AMMI, AIMS Ghana

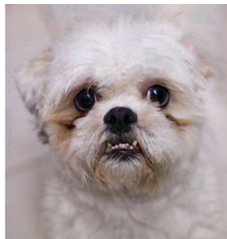
December 6, 2019



dog



+noise



ostrich

# Objectives

1. To show the effect and effectiveness of adversarial examples in machine learning predictions
2. To understand the adversary, and determine how to combat it.
3. To enlighten the audience on machine learning security.

# Outlines

Objectives

Introduction

Properties of Counterfactual Instance

Examples

- Techniques

- Gradient based optimization approach

- Fast gradient sign method

- 1-pixel attack

- Adversarial Patch

- Robust adversarial examples

- Black Box Attacks

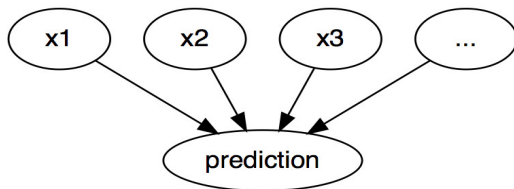
Coding Session

Combating adversarial examples

Conclusion

# Introduction

- An **adversarial example** is an instance with small, intentional feature perturbations that cause a machine learning model to make a false prediction.<sup>1</sup>
- A type of **counterfactual example**



**Figure:** Causal relationships between inputs of a machine learning model and the predictions

---

<sup>1</sup><https://christophm.github.io/interpretable-ml-book/adversarial.html>

# Properties of Counterfactual Instance

A counterfactual should;

- be as **similar** as possible to the instance regarding feature values

# Properties of Counterfactual Instance

A counterfactual should;

- be as **similar** as possible to the instance regarding feature values
- change as **few** features as possible.

# Properties of Counterfactual Instance

A counterfactual should;

- be as **similar** as possible to the instance regarding feature values
- change as **few** features as possible.
- have feature values that are **likely**.

# Properties of Counterfactual Instance

A counterfactual should;

- be as **similar** as possible to the instance regarding feature values
- change as **few** features as possible.
- have feature values that are **likely**.
- produce the predefined prediction as **closely** as possible.



# Examples

1. You submit your details for an offer in such a way that the machine classify you as eligible.

# Examples

1. You submit your details for an offer in such a way that the machine classify you as eligible.
2. A spam detector by-passed

# Examples

1. You submit your details for an offer in such a way that the machine classify you as eligible.
2. A spam detector by-passed
3. Object counterfeit - knife as umbrella

# Examples

1. You submit your details for an offer in such a way that the machine classify you as eligible.
2. A spam detector by-passed
3. Object counterfeit - knife as umbrella
4. Self-driving cars can be deceived by images to misclassify stop-signs.

# Techniques

1. Minimize a distance between the adversarial example generated and the instance to be manipulated

# Techniques

1. Minimize a distance between the adversarial example generated and the instance to be manipulated
2. Perturb the example using the gradients of the model,

# Techniques

1. Minimize a distance between the adversarial example generated and the instance to be manipulated
2. Perturb the example using the gradients of the model,
3. Use the prediction function to train a model to generate new examples,

# Techniques

1. Minimize a distance between the adversarial example generated and the instance to be manipulated
2. Perturb the example using the gradients of the model,
3. Use the prediction function to train a model to generate new examples,

Our focus will be on how adversarial examples affect image classifiers with deep neural networks.



# Gradient based optimization approach

$$\min \text{loss}(f(x + p), y_{adv}) + c \cdot |p|$$

where  $x$  is an image,  $p$  is the changes to the pixels to create an adversarial image,  $y_{adv}$  is the desired outcome class, and the parameter  $c$  is a balancing factor.

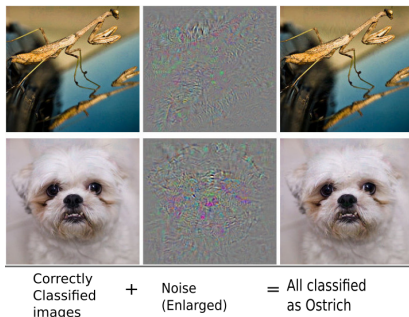
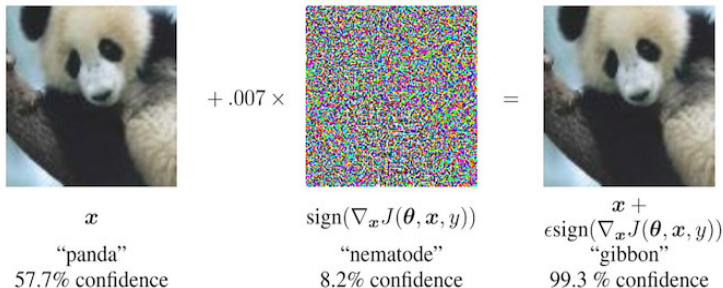


Figure: Examples generated on Alexnet using GB<sup>2</sup>

# Fast gradient sign method

$$x_{adv} = x + \epsilon \text{Sign}(\nabla_x J(\theta, x, y))$$

where  $x$  is the gradient of the models loss function with respect to the original input pixel vector  $x$ ,  $y$  is the true label vector for  $x$  and  $\theta$  is the model parameter vector.



**Figure:** NN predicts Gibbon for a perturbed panda image<sup>3</sup>

<sup>3</sup>Goodfellow et al. "Explaining and harnessing adversarial examples."(2014)

# Changing a single pixel

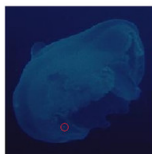
Uses **differential evolution** to find out which pixel is to be changed and how.



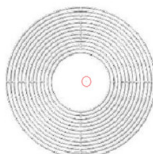
Planetarium  
Mosque(7.81%)



Comforter  
Pillow(6.83%)



Jellyfish  
Bathing tub(21.18%)

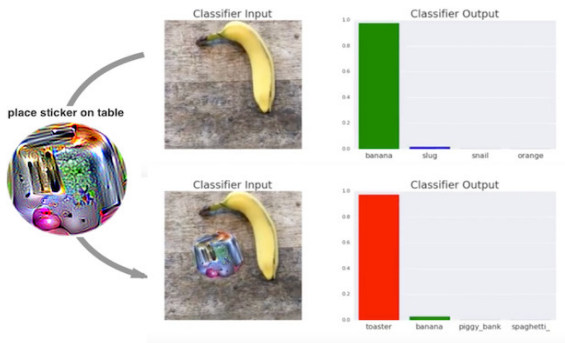


Whorl  
Blower (37.00%)

**Figure:** Changing a single pixel (marked with circles) to deceive a NN to predict the wrong class instead of the original class.<sup>5</sup>

# Adversarial Patch

Replaces a part of the image with a patch that can take on any shape.



**Figure:** Changing a single pixel (marked with circles) to deceive an NN to predict the wrong class instead of the original class.<sup>7</sup>

<sup>7</sup> Brown, Tom B., et al. "Adversarial patch.(2017)

# Robust adversarial examples

- Adversarial over transformations (rotation, zoom in) unlike other methods such as FGM.
- Expectation Over Transformation (EOT) algorithm.



■ classified as turtle      ■ classified as rifle  
■ classified as other

**Figure:** 3D-printed turtle that was designed to look like a rifle to a deep NN<sup>9</sup>

# Black Box Attacks

- No internal model information required and no access to the training data.

# Black Box Attacks

- No internal model information required and no access to the training data.
- Zero access to model gradient

# Black Box Attacks


- No internal model information required and no access to the training data.
- Zero access to model gradient
- A surrogate model is trained to approximate the decision boundaries of the black box model,



# Black Box Attacks

- No internal model information required and no access to the training data.
- Zero access to model gradient
- A surrogate model is trained to approximate the decision boundaries of the black box model,
- Can be used to attack machine learning models on cloud platforms with open api access<sup>10</sup>


---

<sup>10</sup> Papernot, Nicolas, et al. "Practical black-box attacks against machine learning." (2017) 

# Black Box Attacks

- No internal model information required and no access to the training data.
- Zero access to model gradient
- A surrogate model is trained to approximate the decision boundaries of the black box model,
- Can be used to attack machine learning models on cloud platforms with open api access<sup>10</sup>
- Although, Knowledge of domain of input is required


---

<sup>10</sup>Papernot, Nicolas, et al. "Practical black-box attacks against machine learning." (2017) 

# Black Box Attacks

- No internal model information required and no access to the training data.
- Zero access to model gradient
- A surrogate model is trained to approximate the decision boundaries of the black box model,
- Can be used to attack machine learning models on cloud platforms with open api access<sup>10</sup>
- Although, Knowledge of domain of input is required

---

<sup>10</sup>Papernot, Nicolas, et al. "Practical black-box attacks against machine learning." (2017) 

//just let me code

# Combating adversarial examples

AEs can be Model-agnostic.

Methods used to combat adversarial examples include<sup>11</sup>;

- 1 Adversarial training - iterative retraining of the classifier with adversarial examples

---

<sup>11</sup><https://christophm.github.io/interpretable-ml-book/adversarial.html>

# Combating adversarial examples

AEs can be Model-agnostic.

Methods used to combat adversarial examples include<sup>11</sup>;

- 1 Adversarial training - iterative retraining of the classifier with adversarial examples
- 2 Learning invariant transformations of the features or robust optimization (regularization)

---

<sup>11</sup><https://christophm.github.io/interpretable-ml-book/adversarial.html>

# Combating adversarial examples

AEs can be Model-agnostic.

Methods used to combat adversarial examples include<sup>11</sup>;

- 1 Adversarial training - iterative retraining of the classifier with adversarial examples
- 2 Learning invariant transformations of the features or robust optimization (regularization)
- 3 Use of multiple classifiers instead of just one and have them vote the prediction (ensemble)

---

<sup>11</sup><https://christophm.github.io/interpretable-ml-book/adversarial.html>

# Combating adversarial examples

AEs can be Model-agnostic.

Methods used to combat adversarial examples include<sup>11</sup>;

- 1 Adversarial training - iterative retraining of the classifier with adversarial examples
- 2 Learning invariant transformations of the features or robust optimization (regularization)
- 3 Use of multiple classifiers instead of just one and have them vote the prediction (ensemble)

Lot's of research ongoing in this field of Adversarial and ML security.

---

<sup>11</sup><https://christophm.github.io/interpretable-ml-book/adversarial.html>



# Conclusion

- The threats of adversarial examples are real and potent.
- These attacks are not limited to computer-vision but span other areas of ML such as NLP, Reinforcement Learning, Speech Recognition e.t.c.
- Increasing development in this field (but with equivalent sophistication in attack methods).

Think of the many different types of spam emails that are constantly evolving (image spam, header masking etc).

# tHANK yOU



for staying awake