



دینه

Adversarial Attack & Defense

Presented by **Mohammad Khalooei**

Under the supervision of

Prof. Mohammad Mehdi Homayounpour & Dr. Maryam Amirmazlaghani



Amirkabir
Artificial Intelligence
Summer
Summit 2019

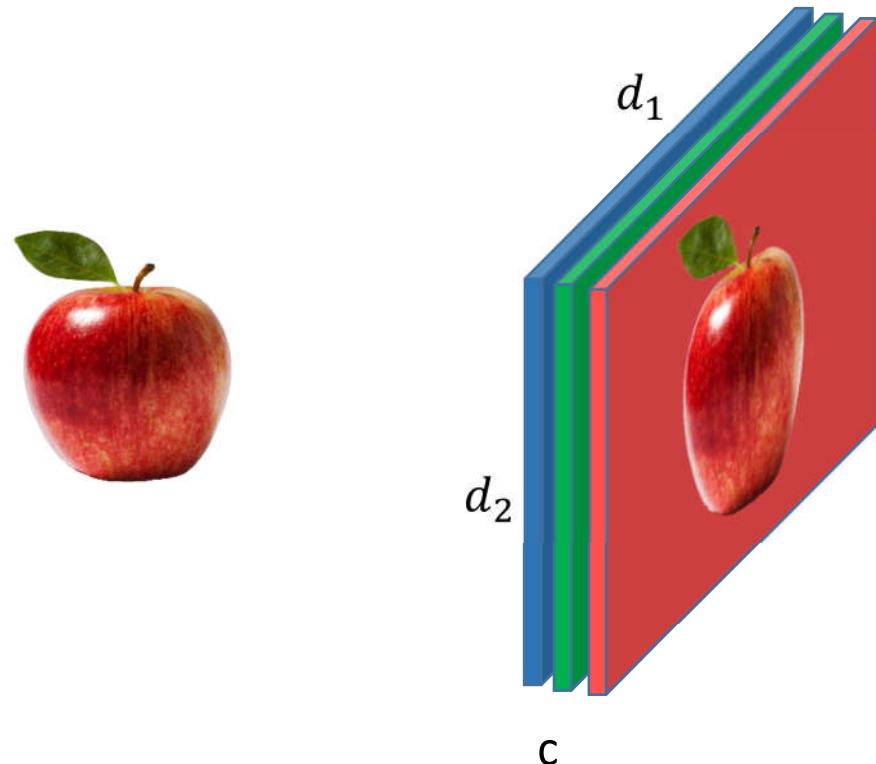
23 - 25 July 2019
Tehran, IRAN

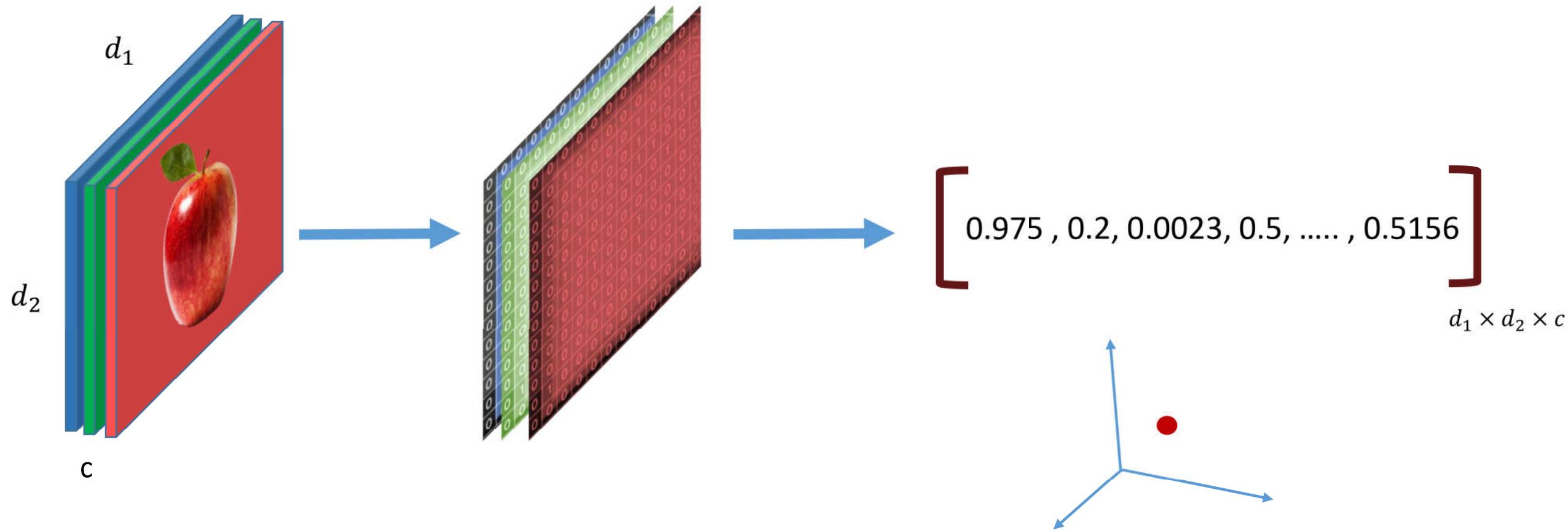




- Overview of Generative & Discriminative models
- Decision-boundary crossing
- Vulnerability of deep neural networks in different data
- Different attacks and penetration ways of ML & DL approaches
- Attack Toolboxes
- Different defense strategies for DNN
- Tips to stay safe in DNN













<https://www.slideshare.net/khalooei/life-of-points-machine-learning-with-orange-flavor>



ظهوری با هوش مصنوعی

زندگی نقطه‌ها

khalooei@aut.ac.ir

<https://ceit.aut.ac.ir/~khalooei>



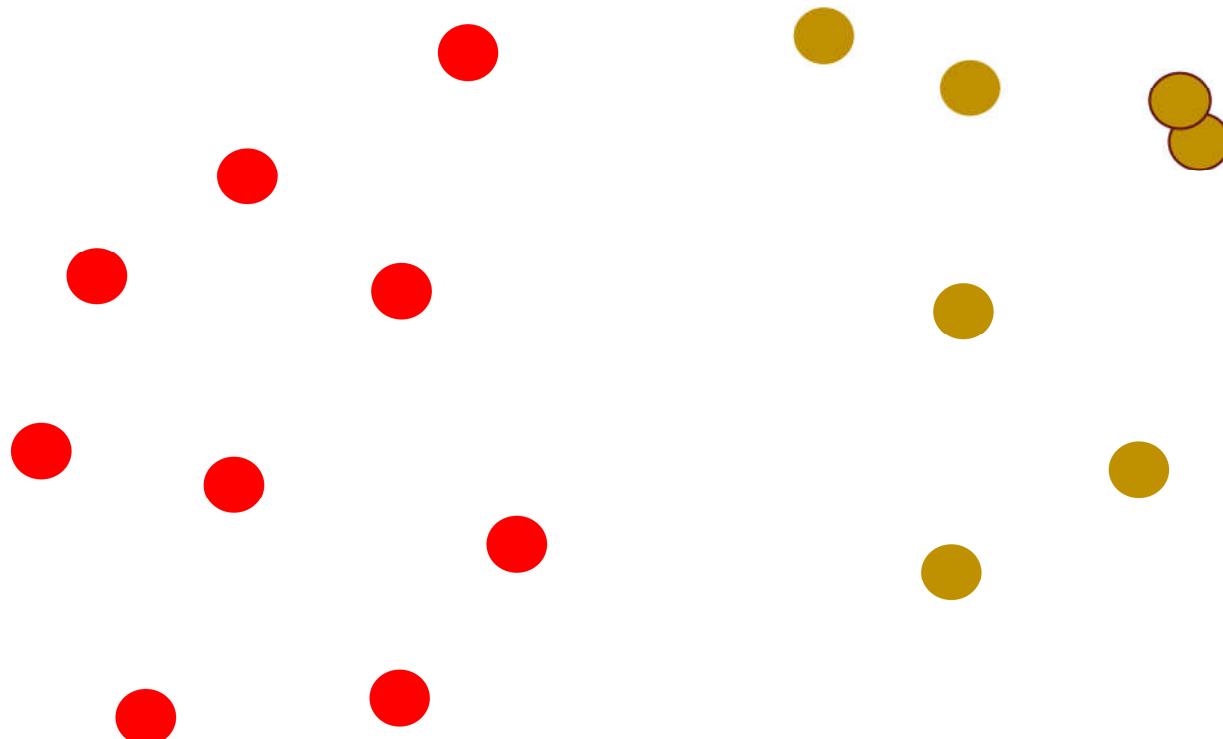
SSC 

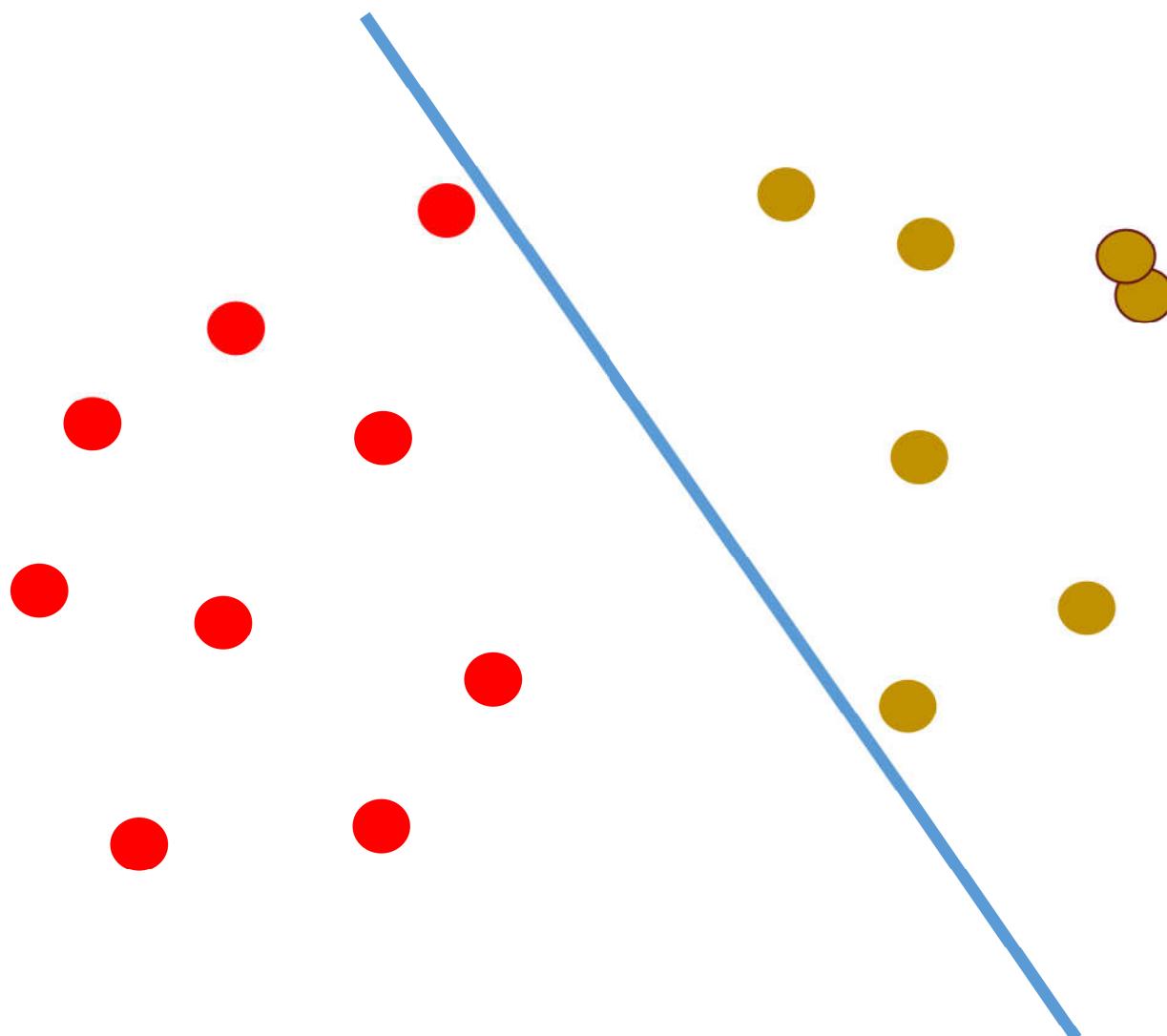


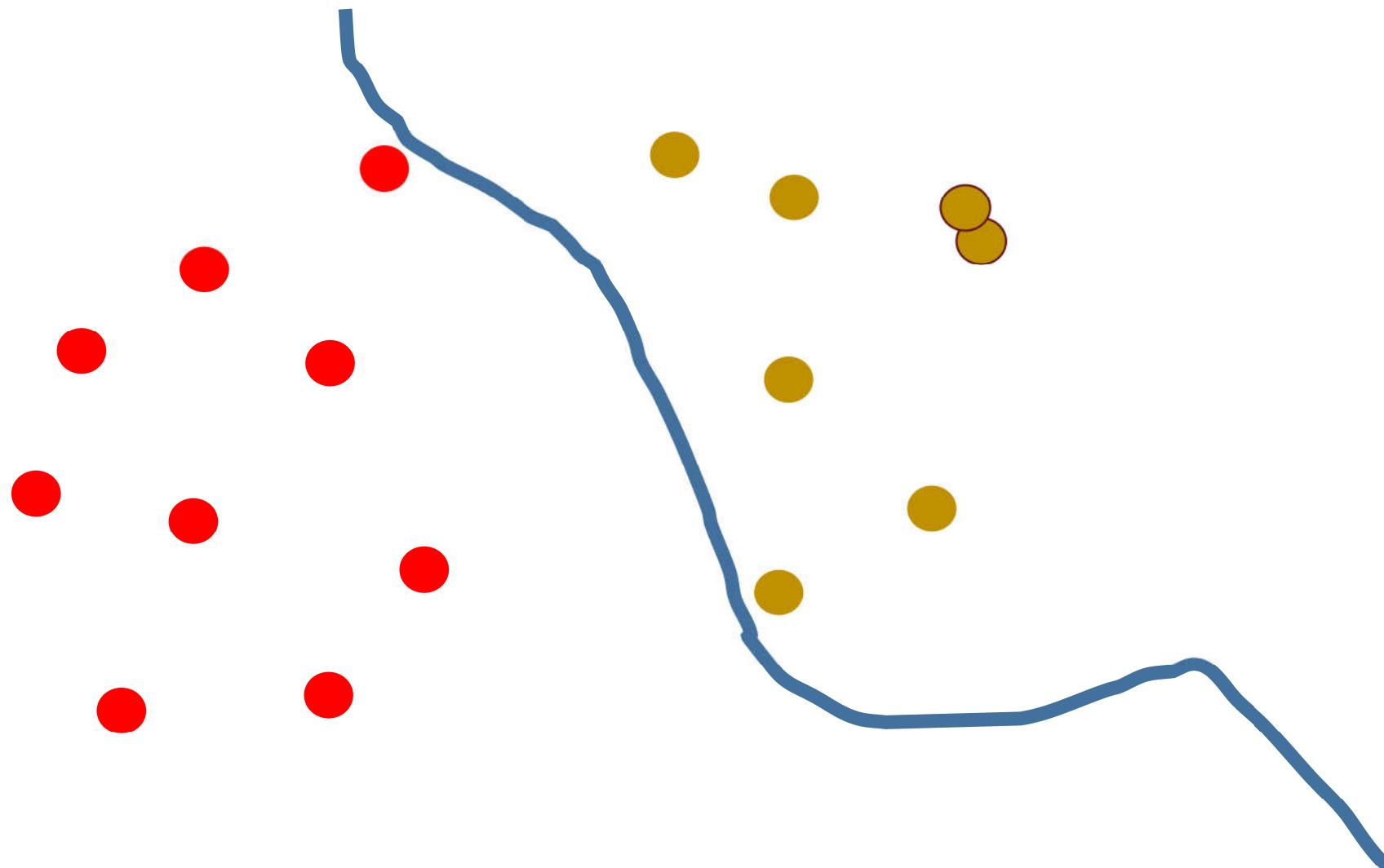


Review:: Discriminative Model



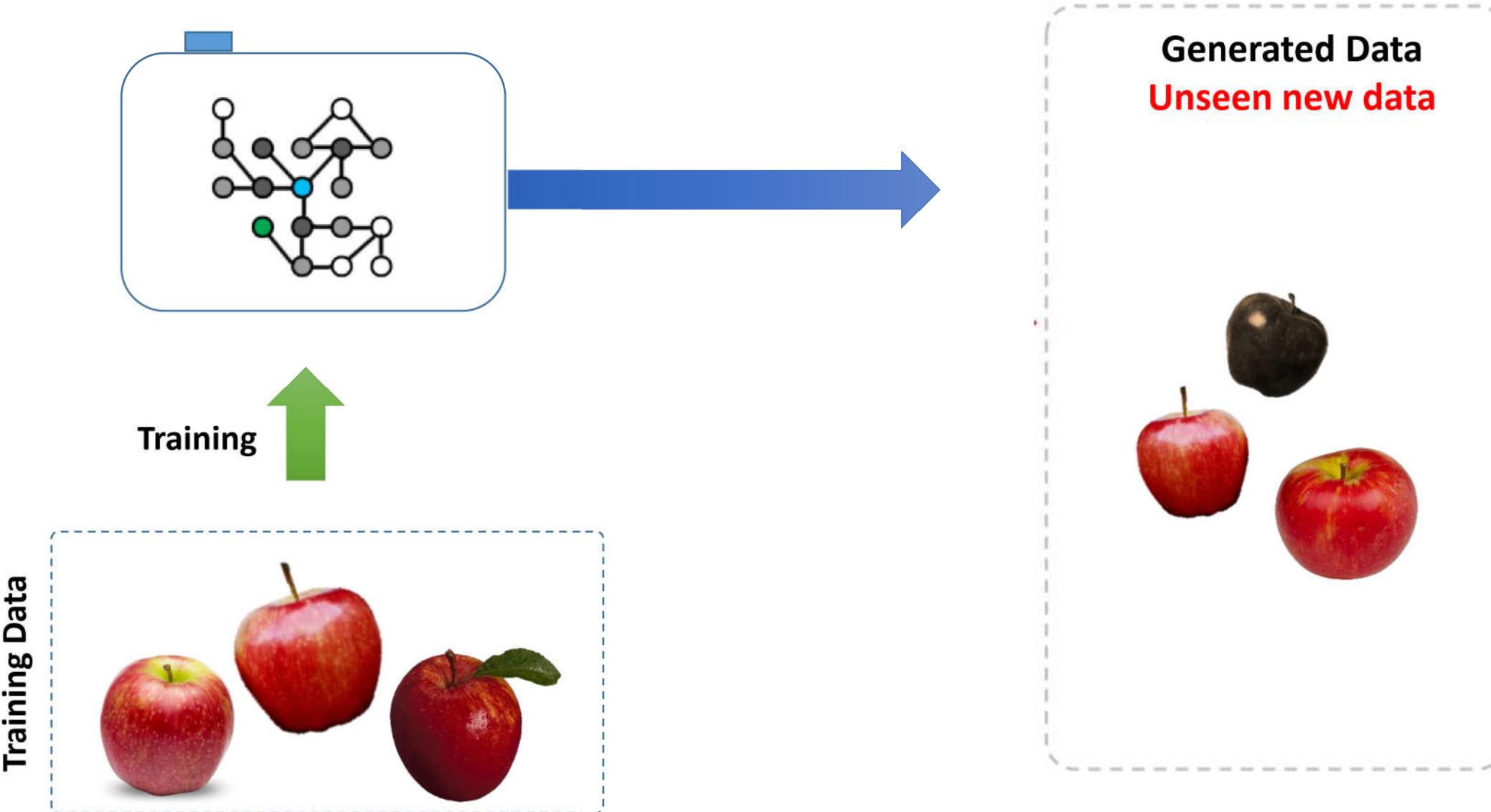




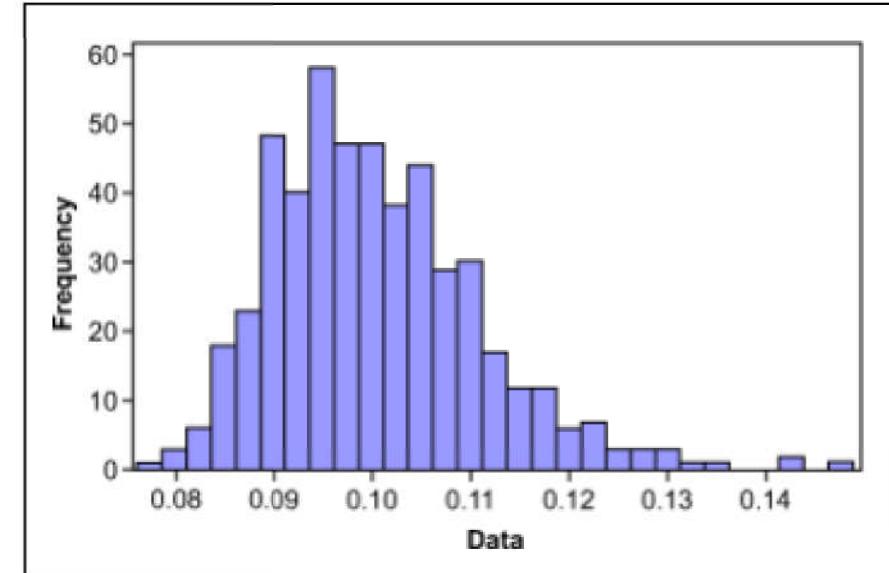




Generative Model

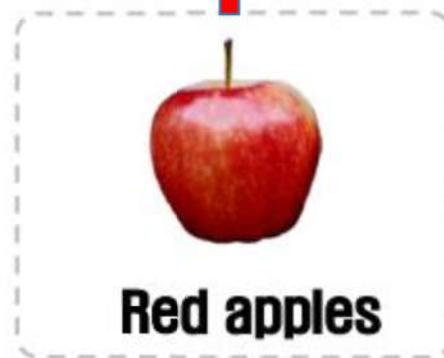
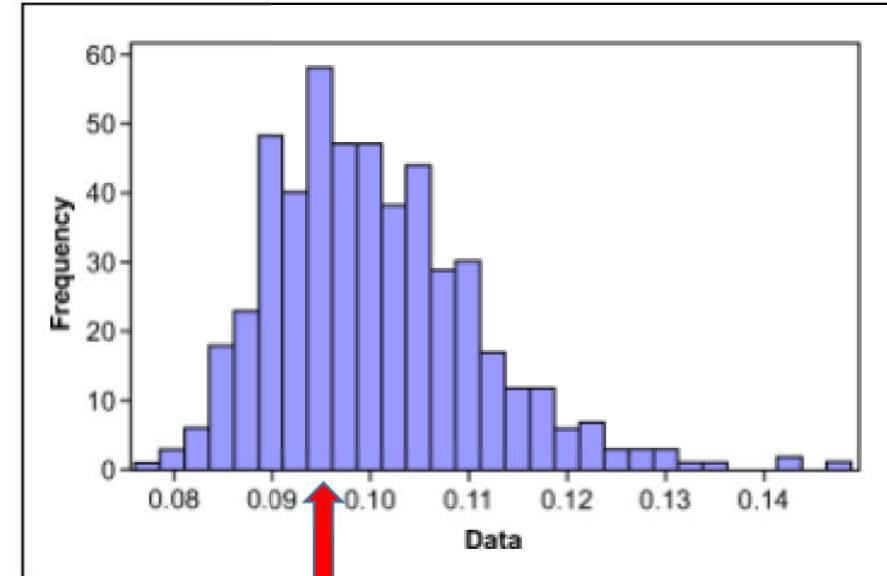
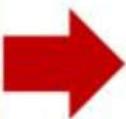


Review:: Generative Model

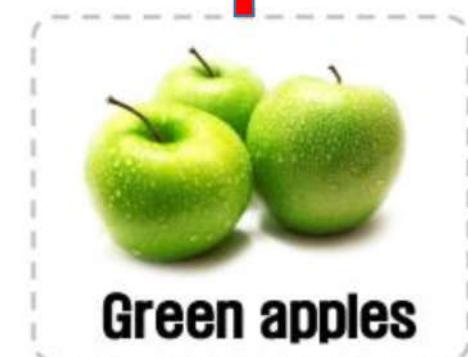
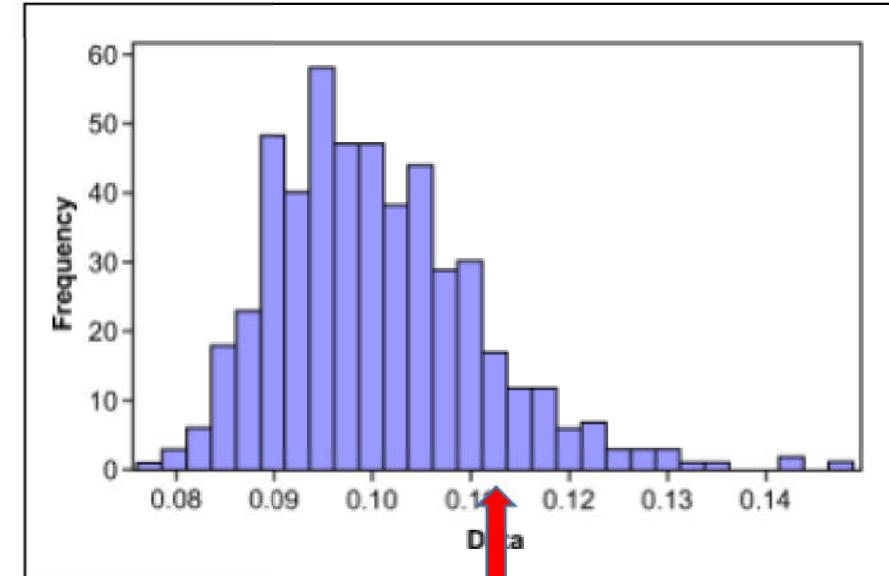


Distribution of the actual images

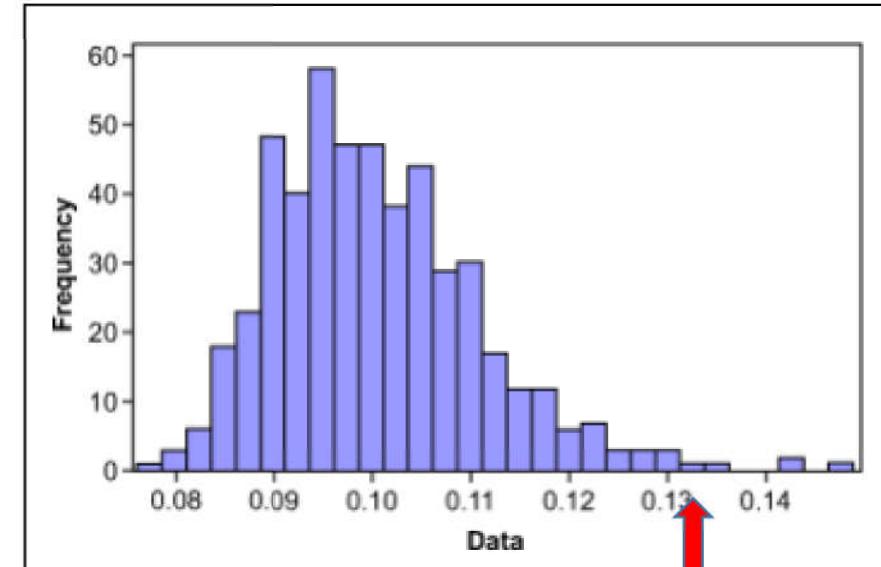
Review:: Generative Model



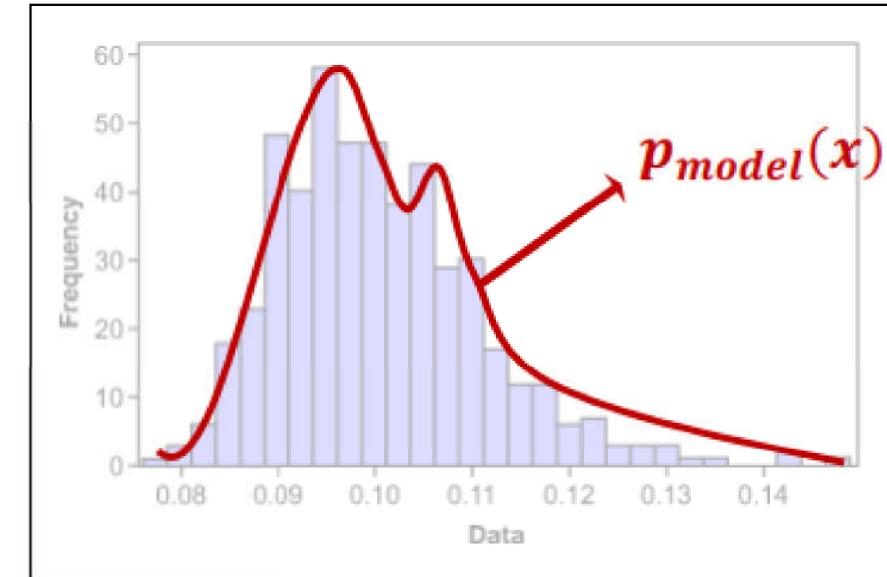
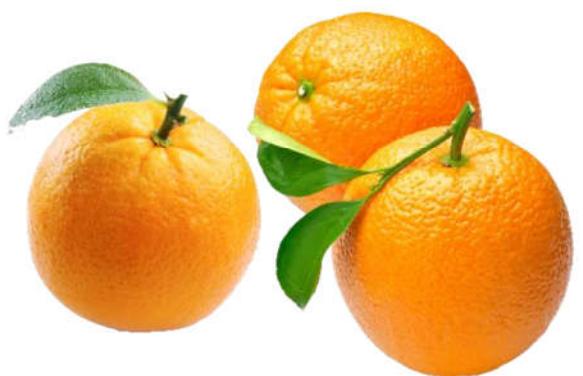
Review:: Generative Model



Review:: Generative Model



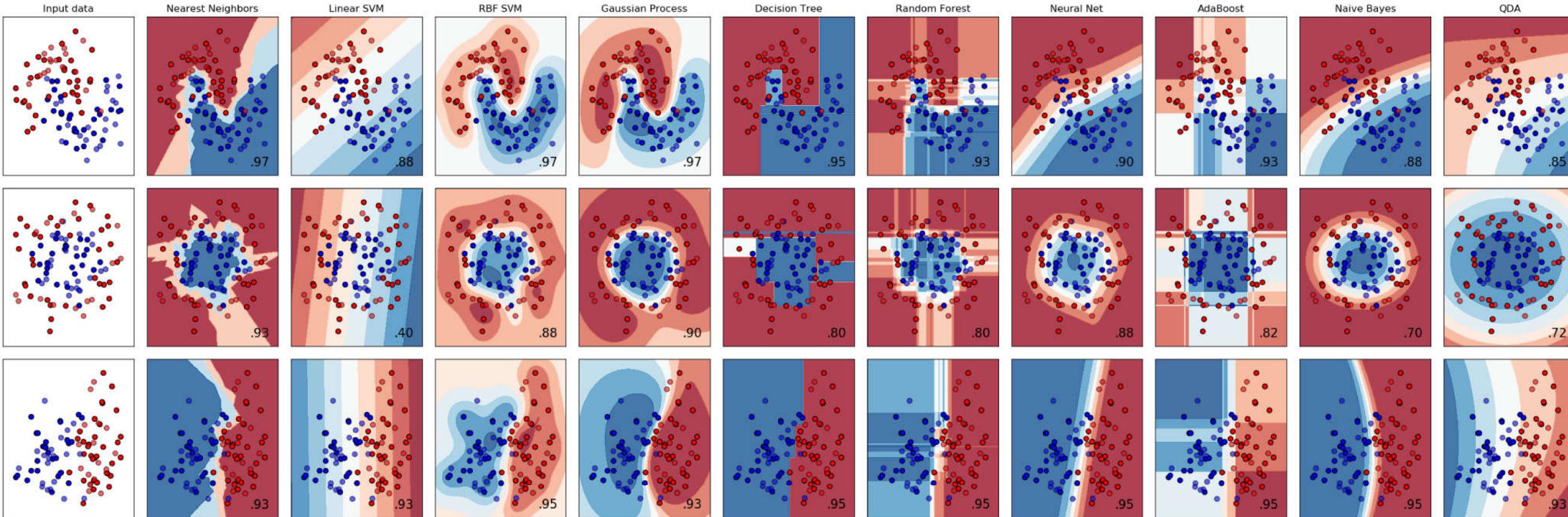
Weird apples



Distribution of the actual images



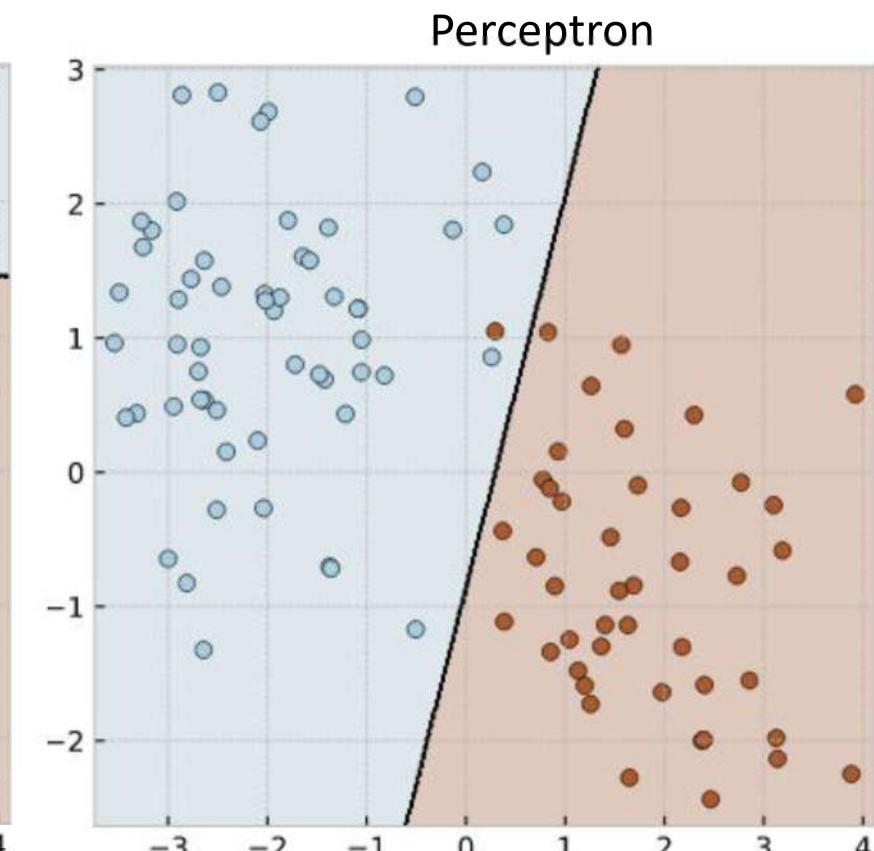
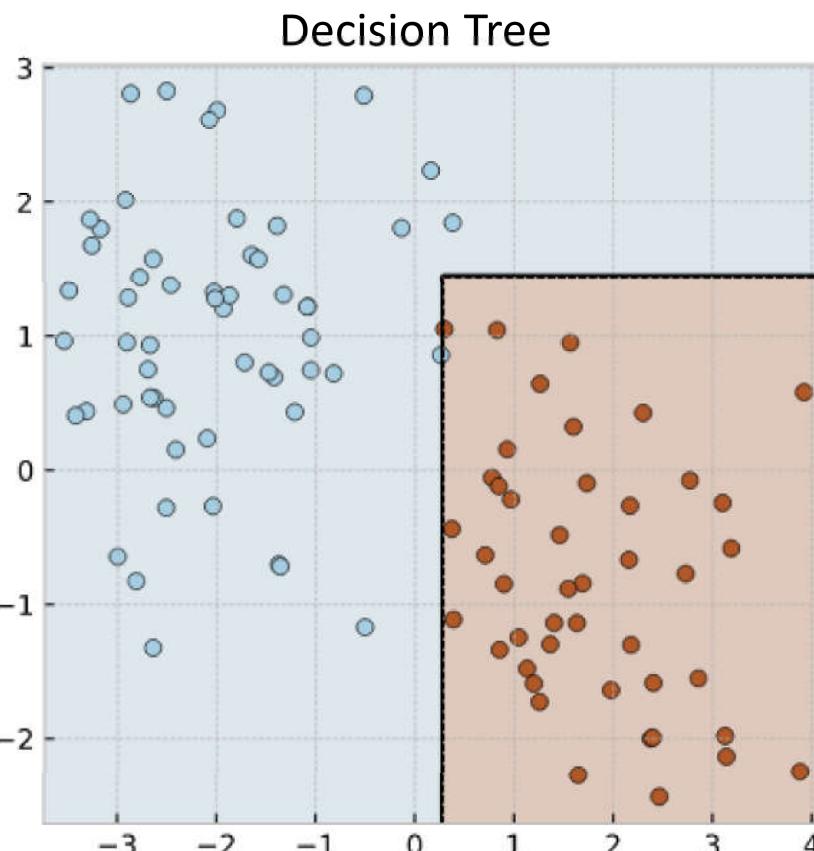
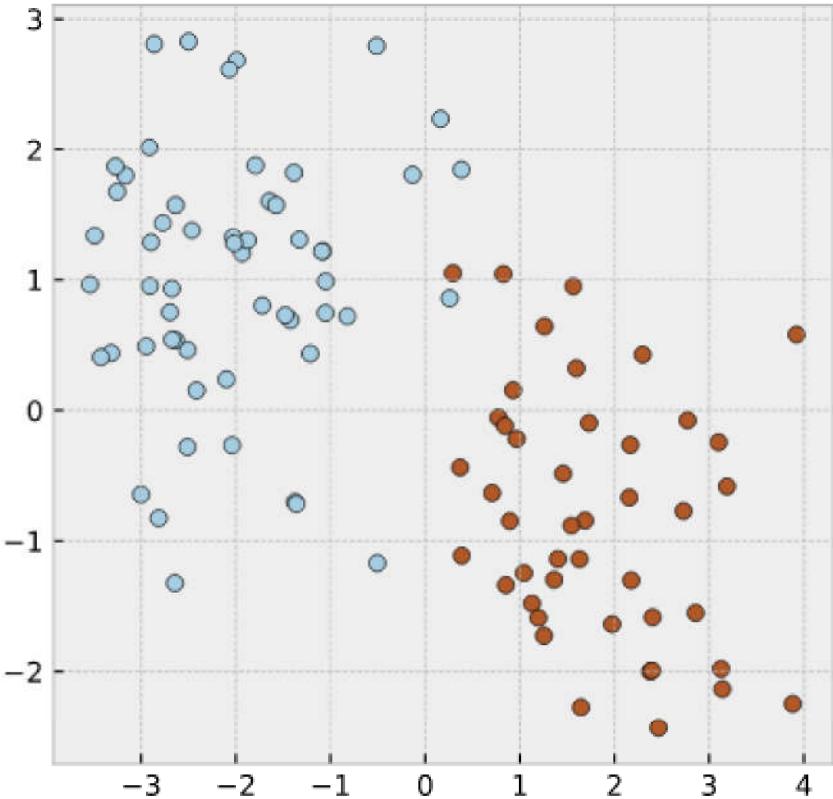
Decision boundary crossing



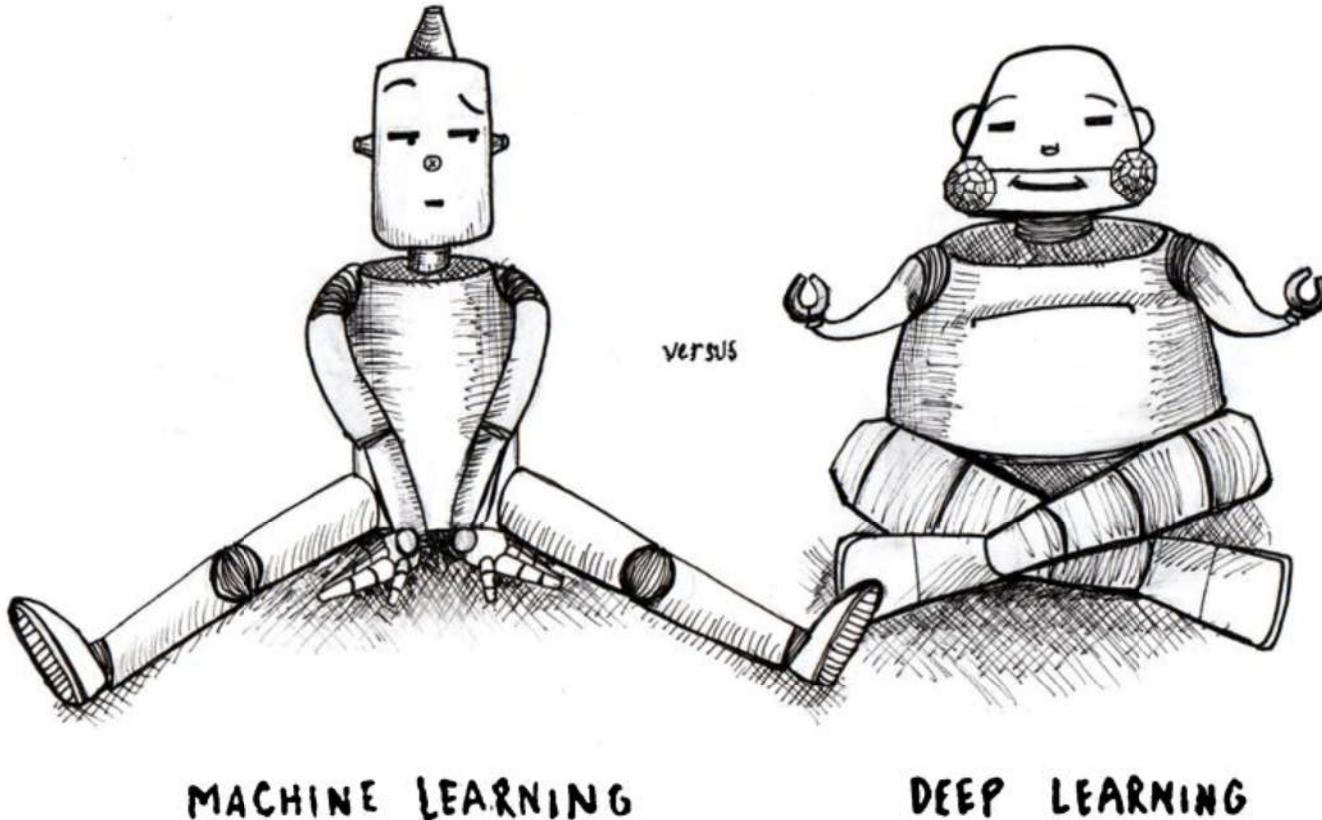
https://scikit-learn.org/stable/auto_examples/classification/plot_classifier_comparison.html



Decision boundary crossing



How can developments in deep learning make for a better approach to value investing?



EUCLIDEAN TECHNOLOGIES MANAGEMENT © 2018



PayPlug

PayPal

 invincea deepinstinct

FeatureSmith

 amazon
web services™

Google Cloud Platform

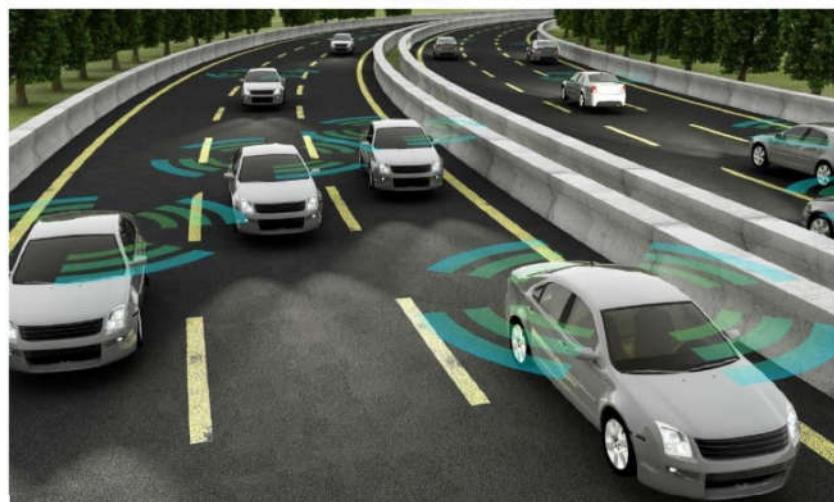
Autonomous
driving

Financial fraud
detection

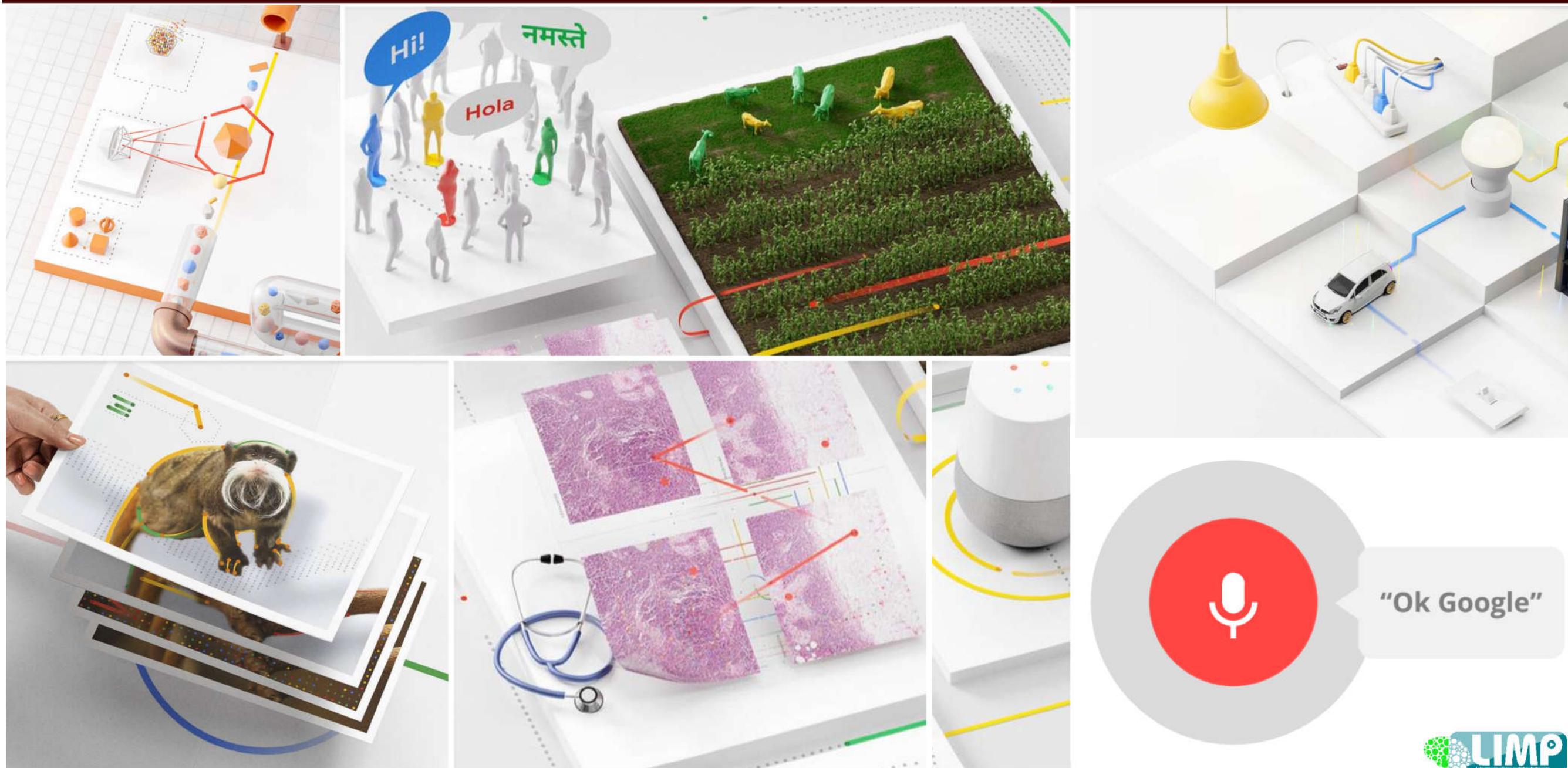
Malware / APT
detection

Machine Learning
as a Service

AI and the top innovation



AI and the top innovation





Adversarial examples represent worst-case domain shifts





Deep learning researchers status! 😊





December 20 '17

Fooling Google's image-recognition AI 1000x faster



Figure 8. The Google Cloud Vision Demo labelling on the unperturbed image.



Skiing	91%
Ski	89%
Piste	86%
Mountain Range	86%
Geological Phenomenon	85%
Glacial Landform	84%
Snow	82%
Winter Sport	78%
Ski Pole	75%

Dog	91%
Dog Like Mammal	87%
Snow	84%
Arctic	70%
Winter	67%
Ice	65%
Fun	60%





November 02 '17

Fooling neural networks w/3D-printed objects

The team 3D-printed a toy turtle that was misclassified as a rifle and a baseball that was classified as an espresso, no matter what angle the neural network views them from.

<https://www.csail.mit.edu/news/fooling-neural-networks-w3d-printed-objects>

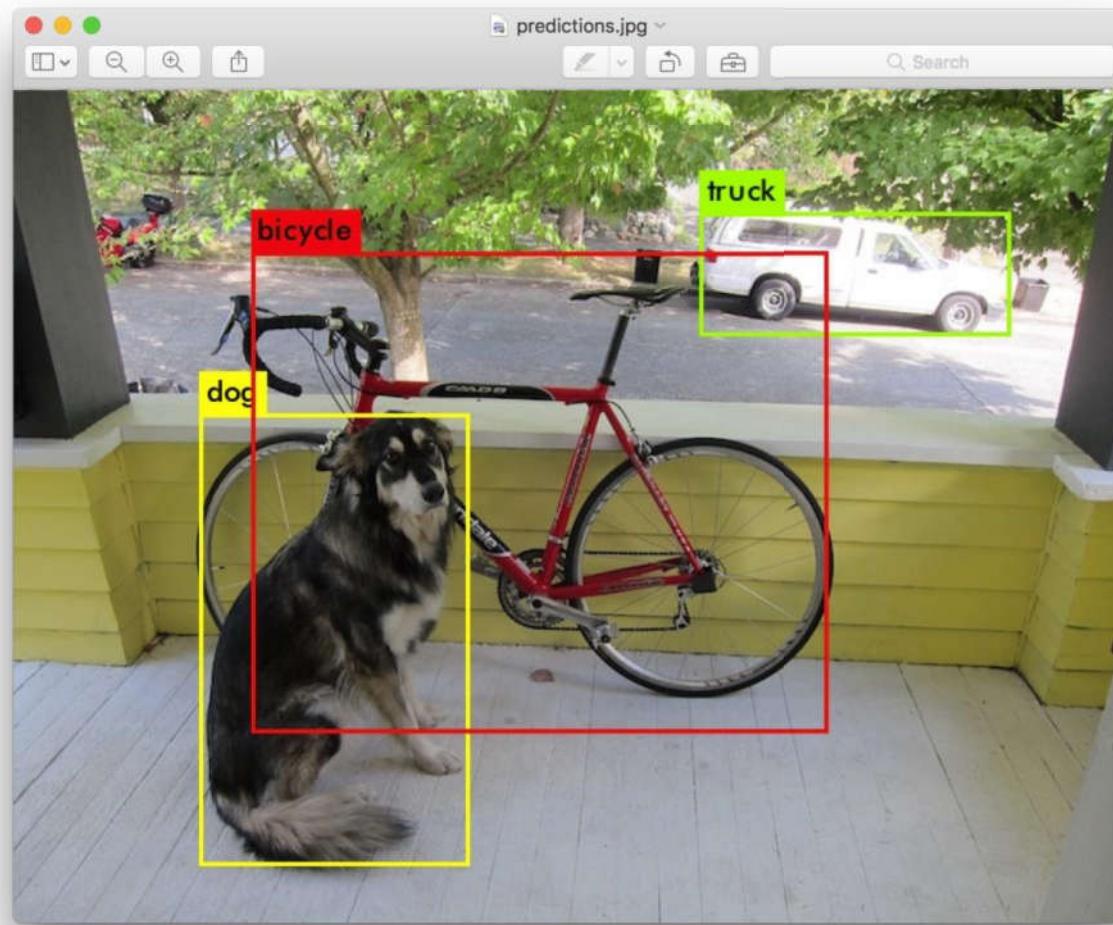
CSAIL :: MIT Computer Science & Artificial Intelligence Lab

- YOLO: Real-Time Object Detection

Redmon J, Divvala S, Girshick R, Farhadi A. You only look once: Unified, real-time object detection. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 2016



<https://pjreddie.com/darknet/yolo/>





Fooling automated surveillance cameras: adversarial patches to attack person detection

Simen Thys*

simen.thys@student.kuleuven.be

Wiebe Van Ranst*

wiebe.vanranst@student.kuleuven.be

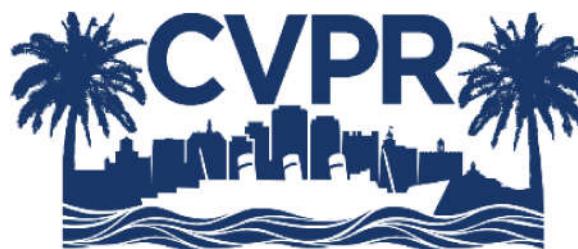
Toon Goedemé

toon.goedeme@student.kuleuven.be

KU Leuven

EAVISE, Technology Campus De Nayer, KU Leuven, Belgium.

* Authors contributed equally to this paper.



LONG BEACH
CALIFORNIA
June 16-20, 2019



<https://syncedreview.com/2019/04/24/now-you-see-me-now-you-dont-fooling-a-person-detector/>





Cloakwear

Home

Catalog



Featured Item



Adversarial T-shirt

From \$30.44



Adversarial Long Sleeve Shirt

From \$34.71



Adversarial Hoodie

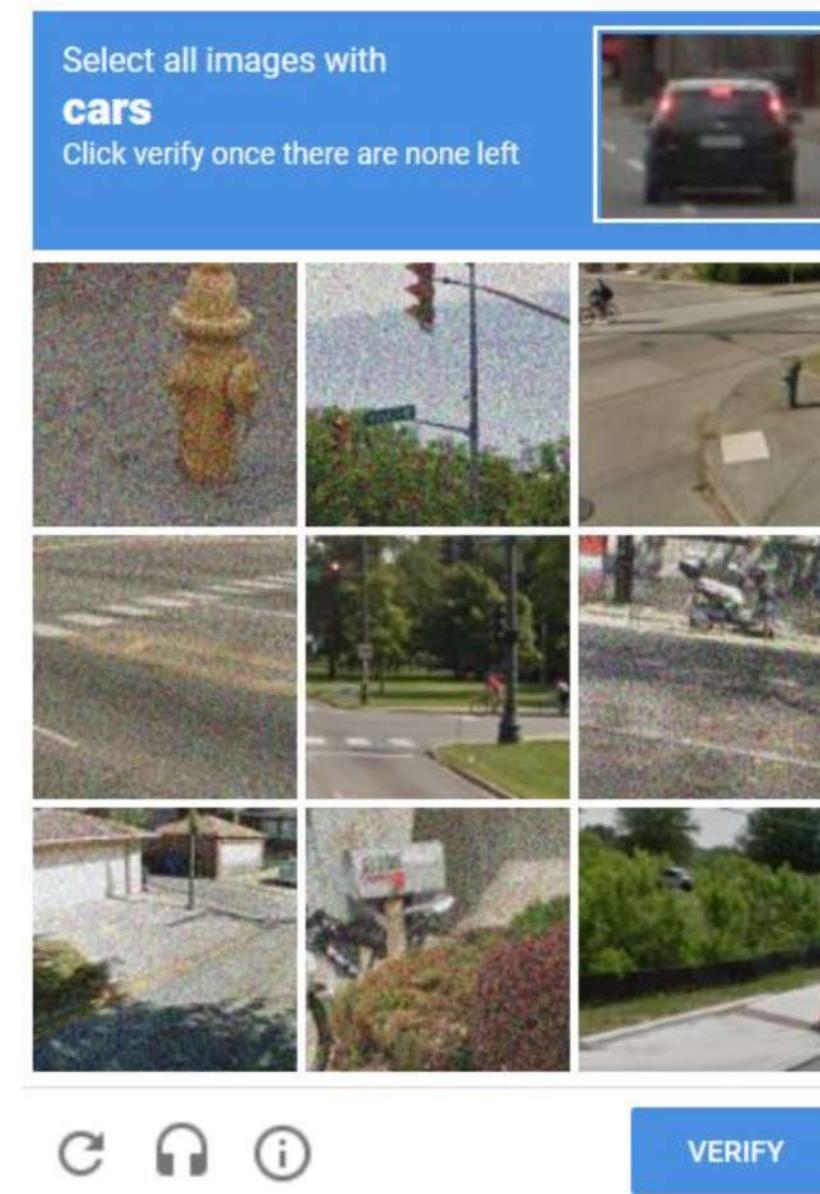
From \$43.39

<https://cloakwear.co/>





An interesting usage :: Google verification!





Deep learning researchers status! 😊

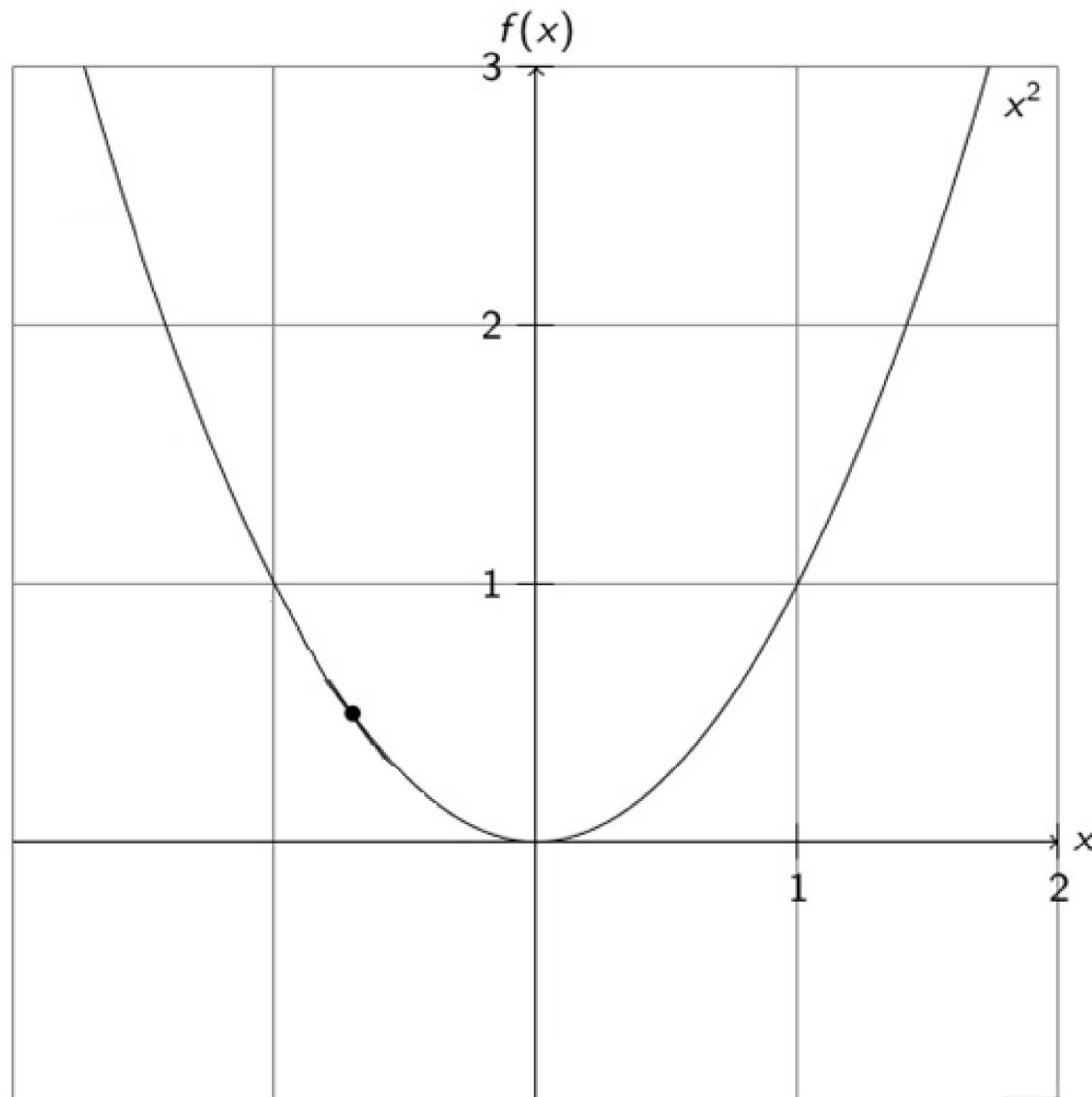


Delving into the problem ...



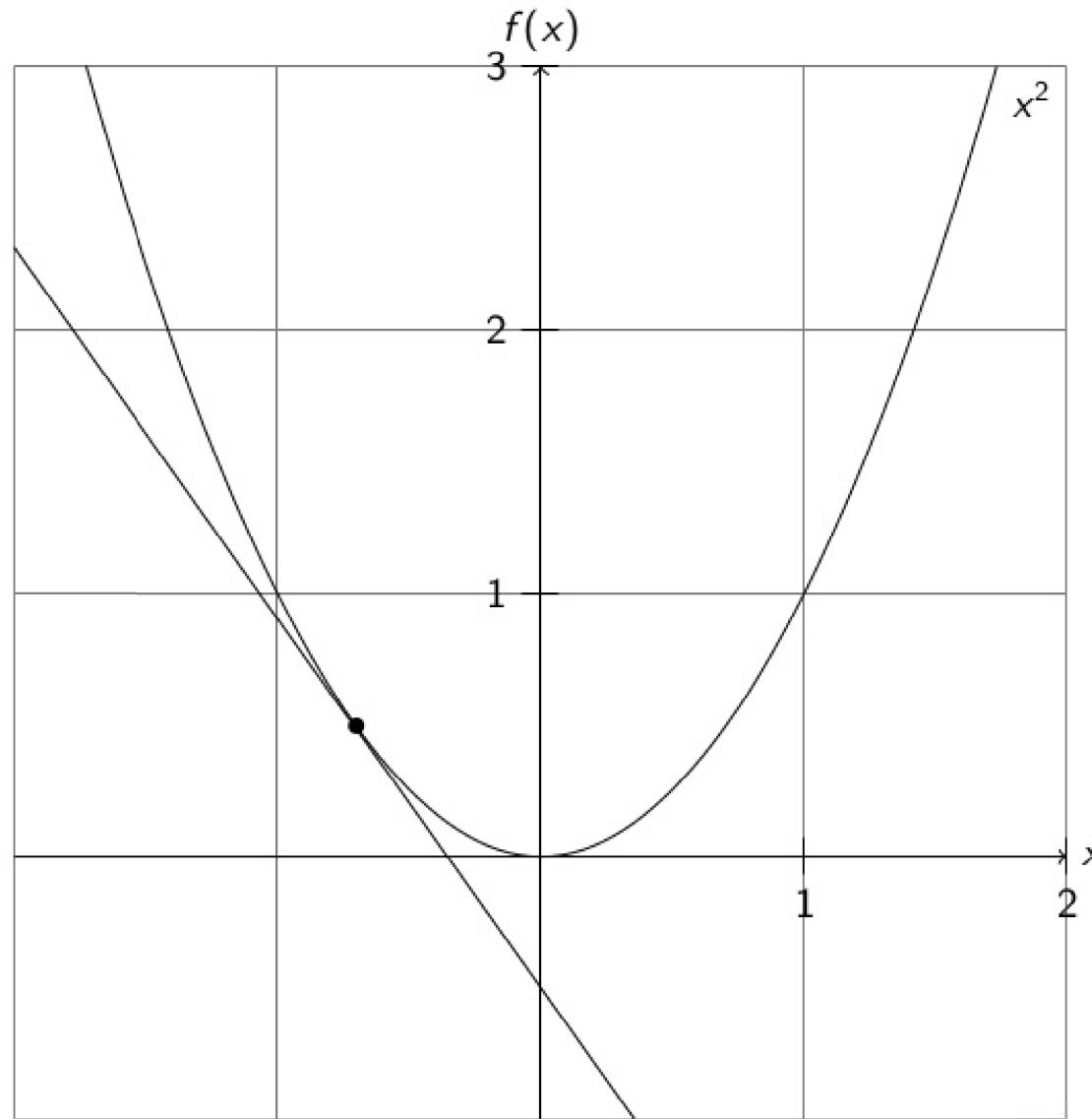


Motivation



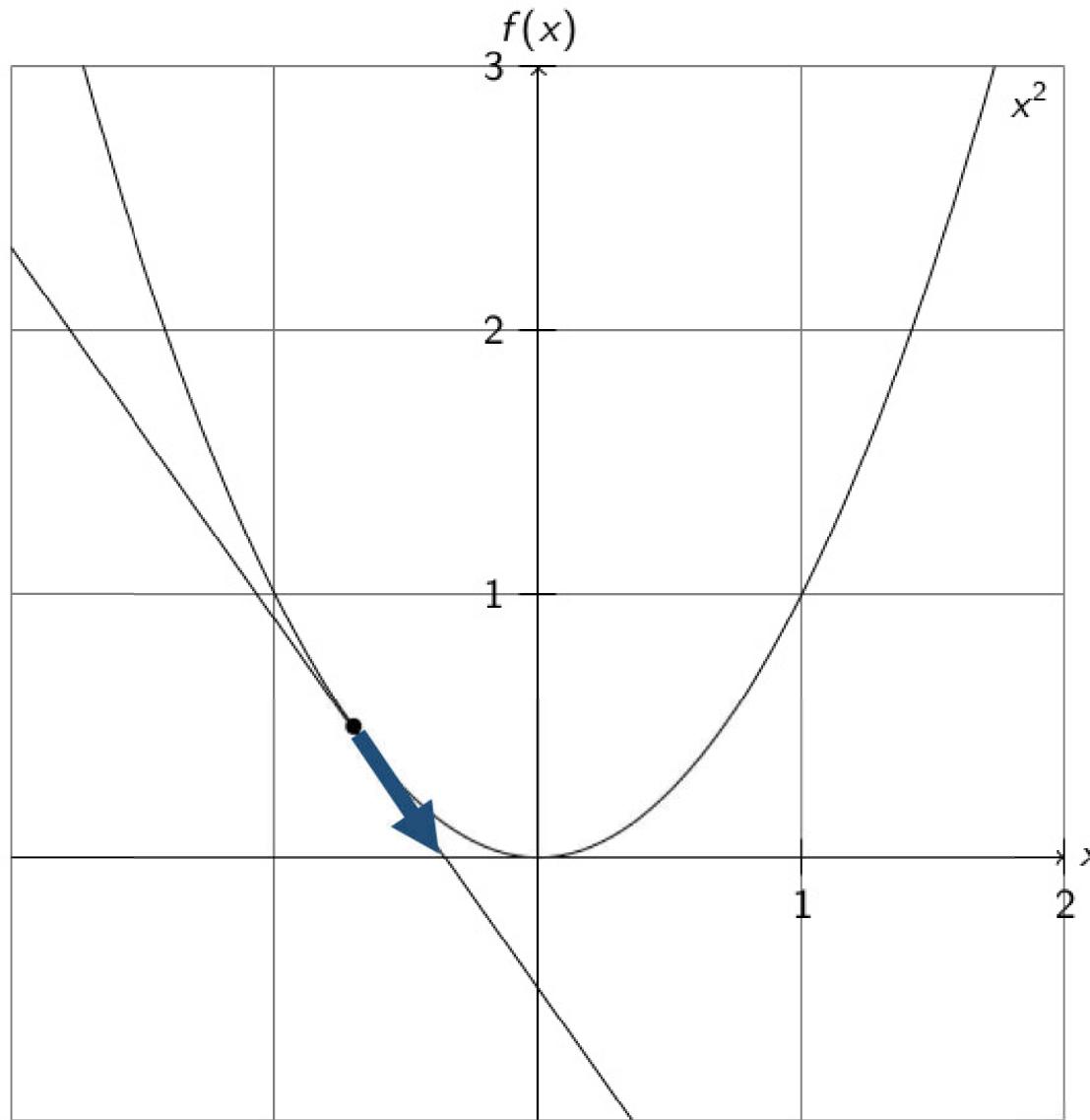


Gradient-descent



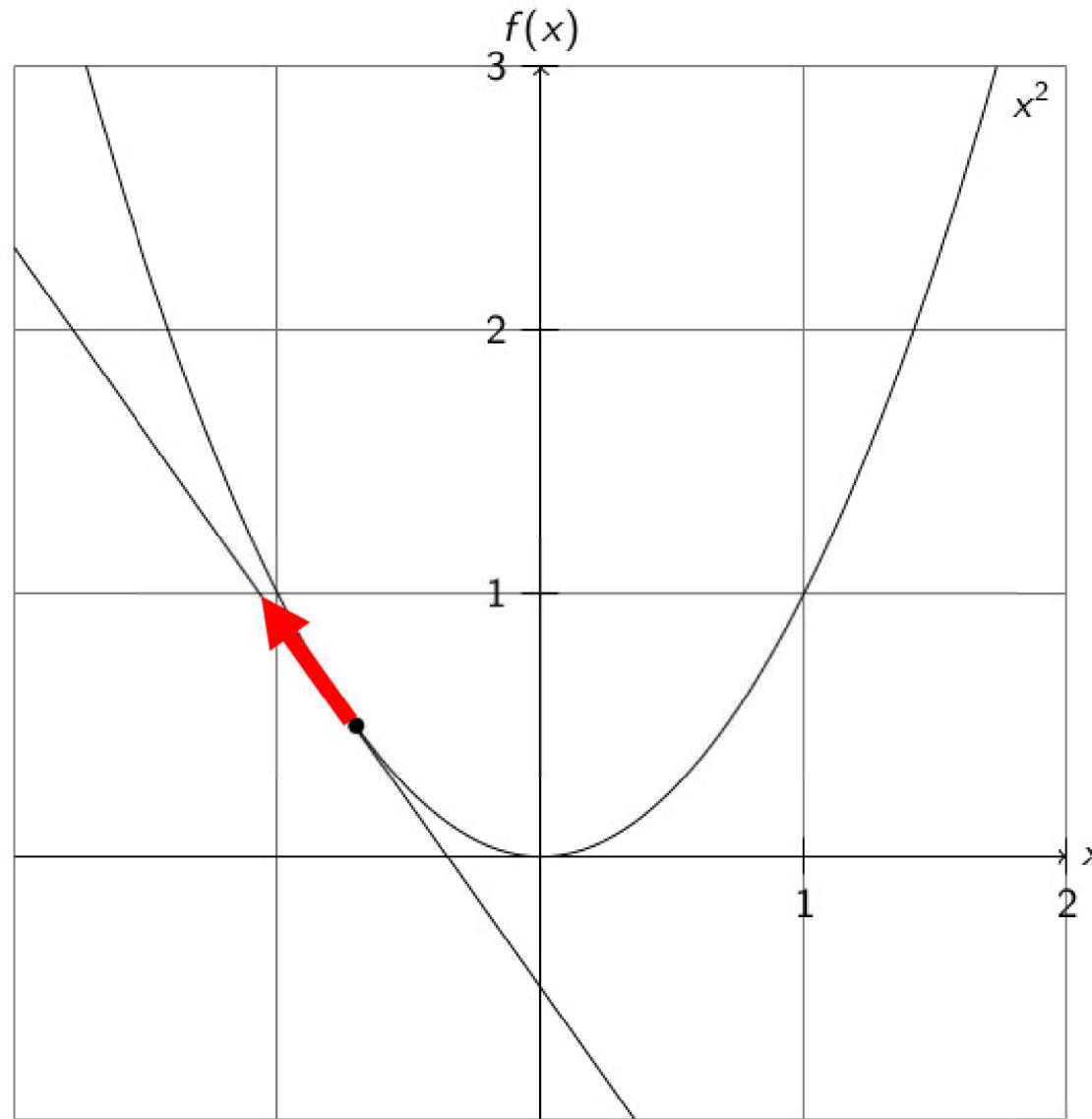


Gradient-descent





Gradient-descent



Attacks ?





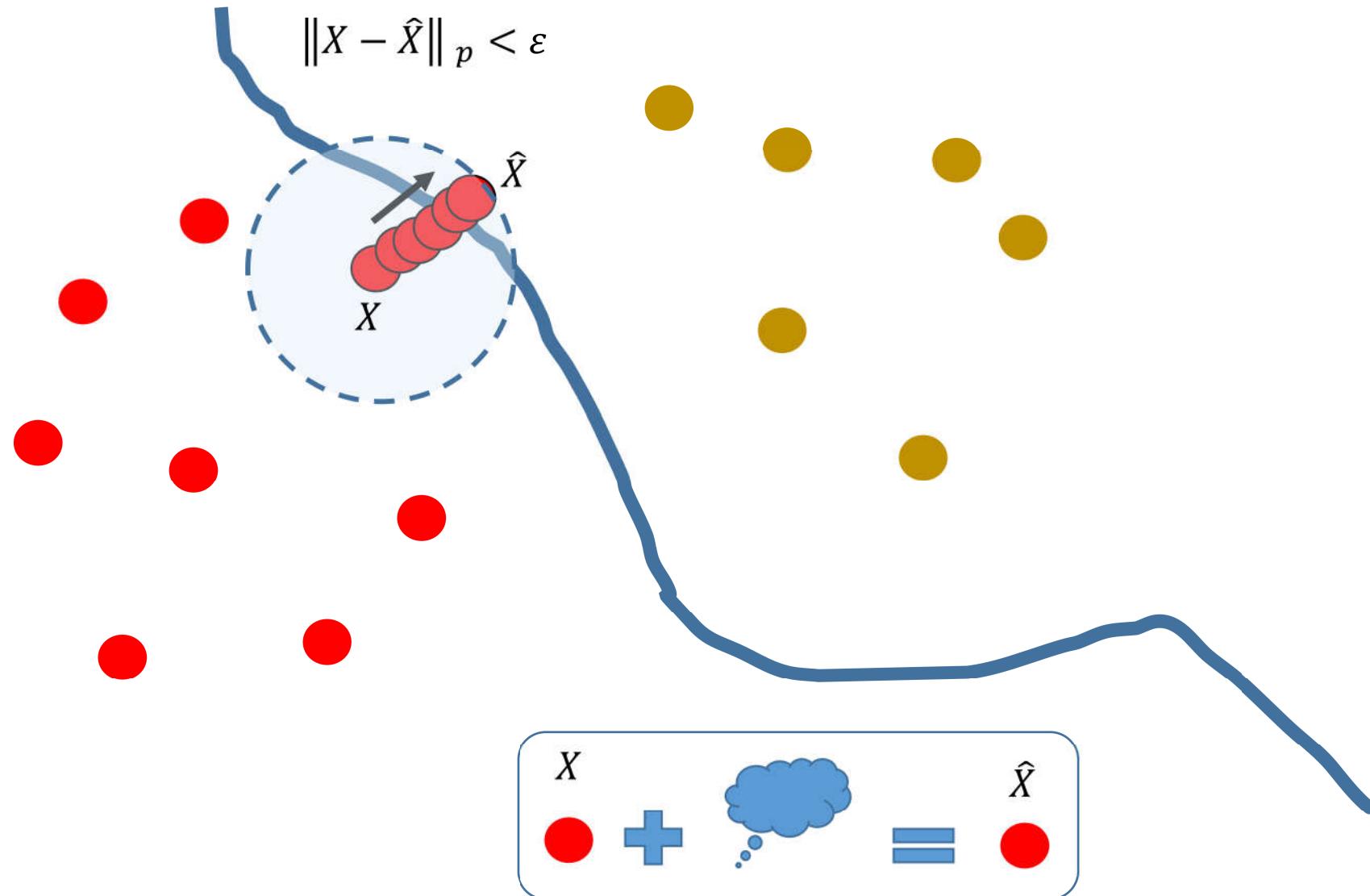
Review the Clever Hans story!



(“Clever Hans,
Clever
Algorithms”, Bob
Sturm)

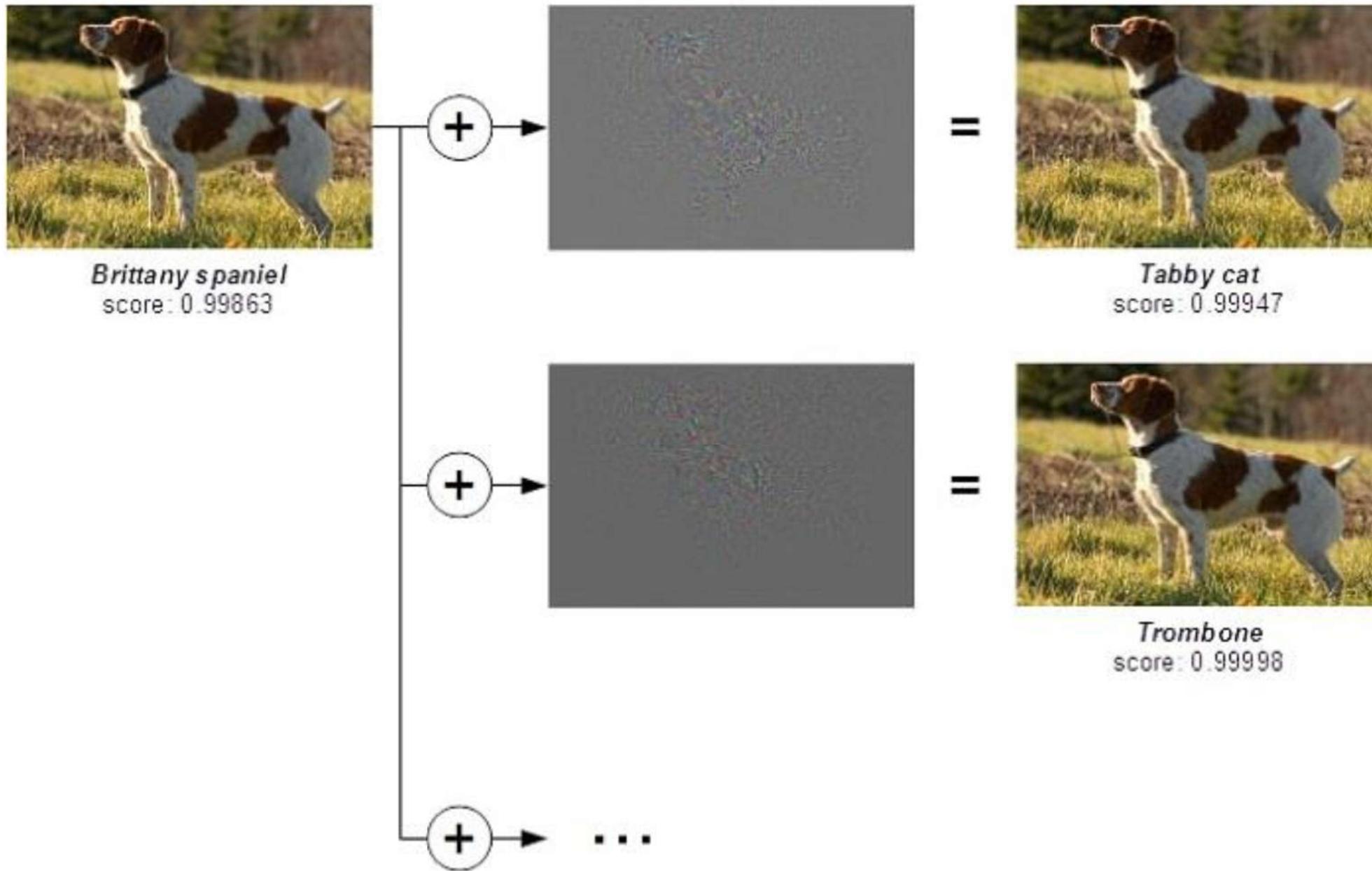


What is Adversarial Example?

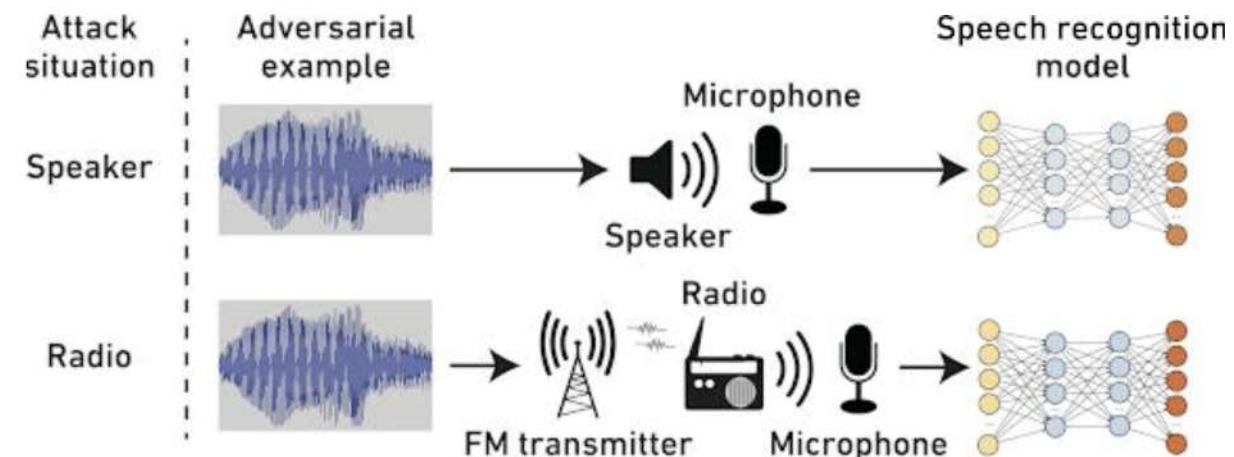




Examples of Adversarial Example



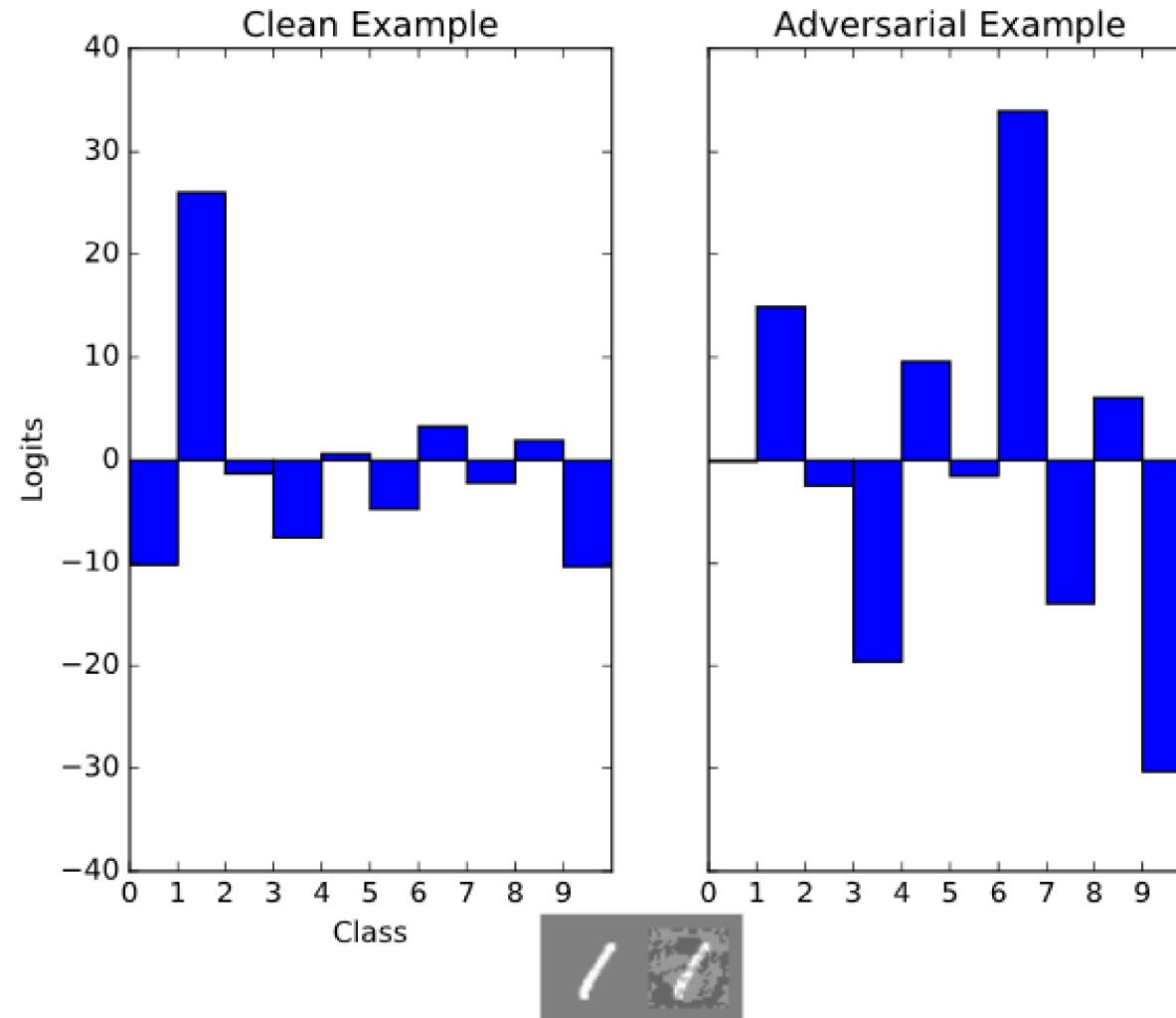
<https://www.kdnuggets.com/2018/10/adversarial-examples-explained.html>



From WWW



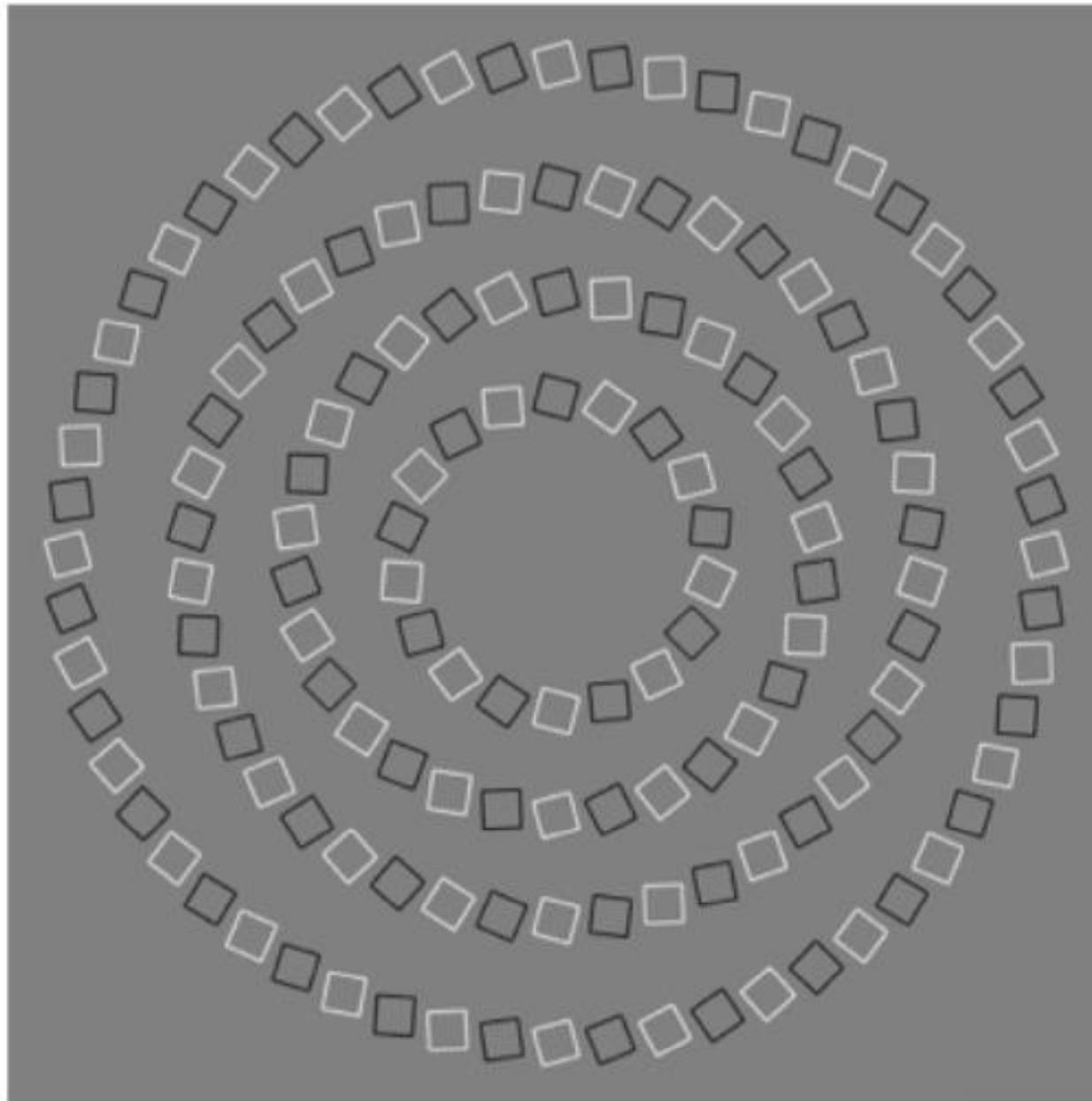
Perturbation's effect on class distributions





Adversarial examples in the human visual system

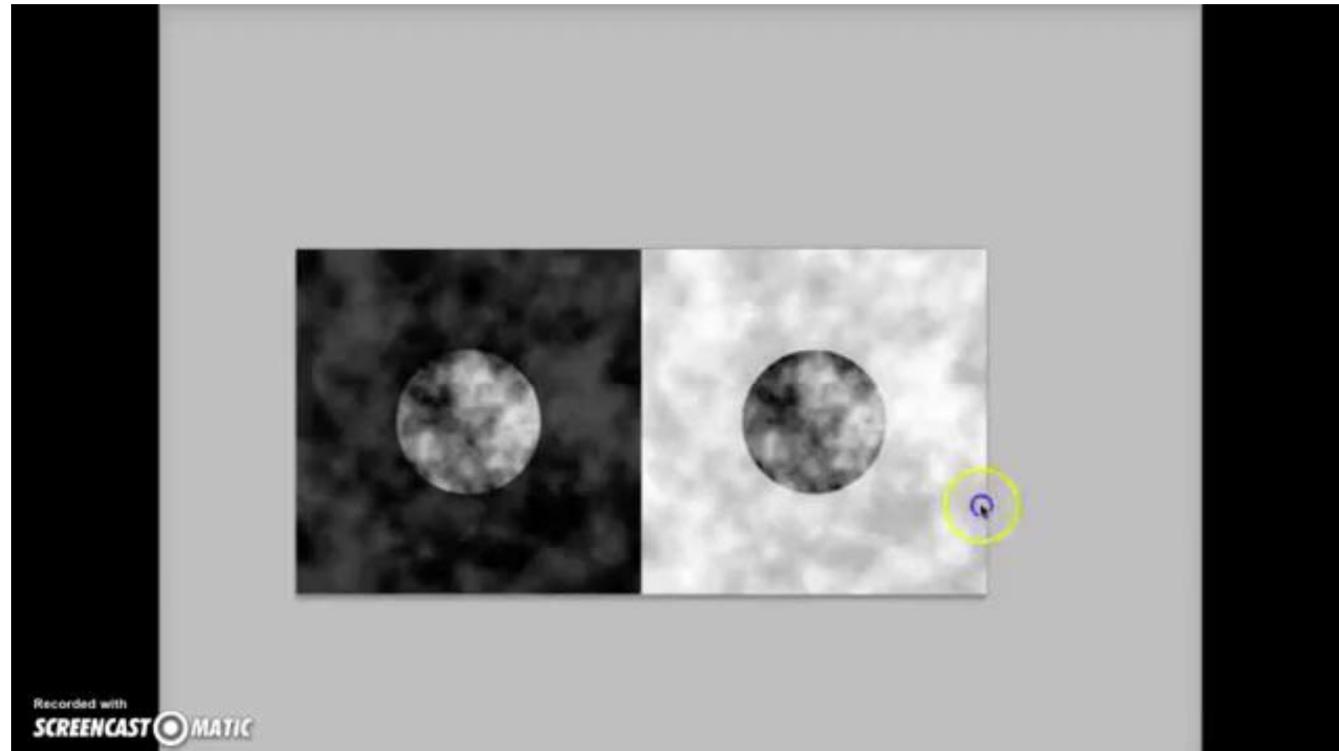
- Circles are concentric but appear intertwining



(Pinna and Gregory, 2002)



Attack: Through the human eye

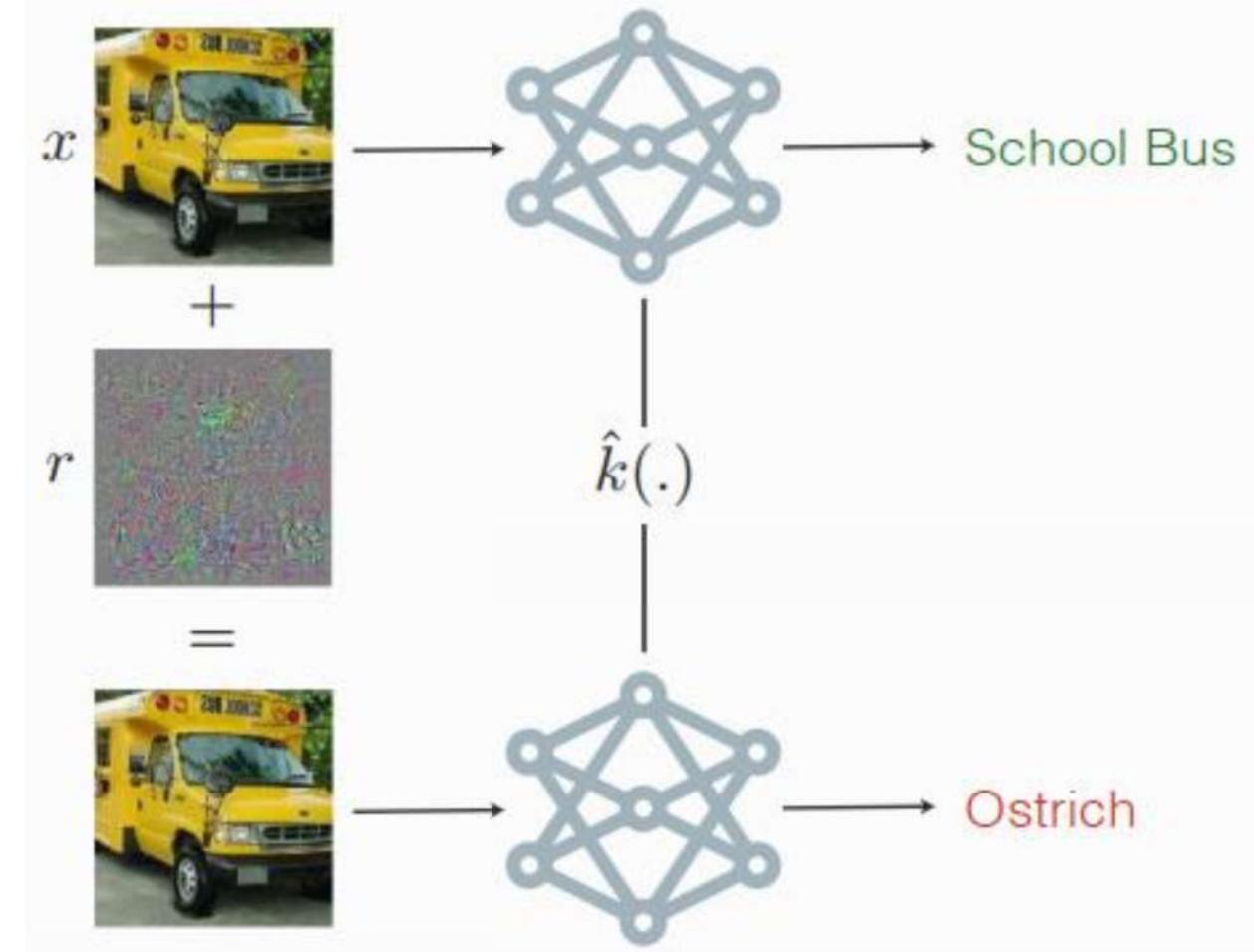


<https://www.slideshare.net/DavidKim486/universal-adversarial-perturbation>



Interiguing properties of Neural Network

C. Szegedy et al. ICLR, 2014



$$r^* = \arg \min_r J(\hat{k}(x + r), y) + C\|r\|$$



Interiguing properties of Neural Network

C. Szegedy et al. ICLR, 2014



Explaining and Harnessing Adversarial Example

I.J. Goodfellow et al. ICLR, 2015

$$\tilde{x} = x + \eta$$

Perturbation

$$\eta = \epsilon \text{sign}(\nabla_x J(\theta, x, y))$$

Gradient of the cost function

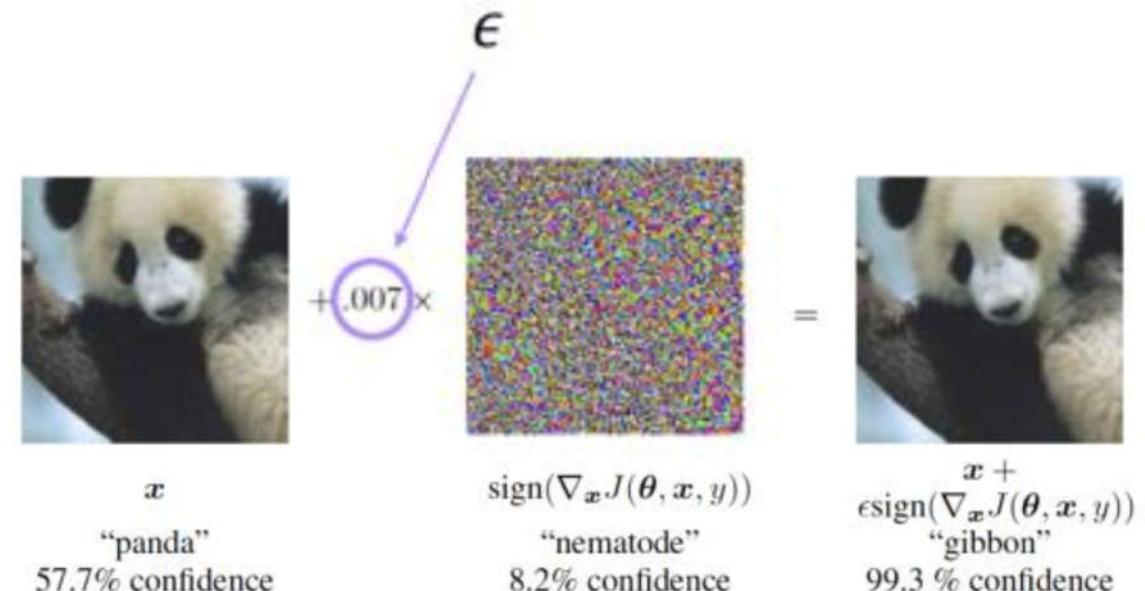


Figure 1: A demonstration of fast adversarial example generation applied to GoogLeNet (Szegedy et al., 2014a) on ImageNet. By adding an imperceptibly small vector whose elements are equal to the sign of the elements of the gradient of the cost function with respect to the input, we can change GoogLeNet's classification of the image. Here our ϵ of .007 corresponds to the magnitude of the smallest bit of an 8 bit image encoding after GoogLeNet's conversion to real numbers.



Interiguing properties of Neural Network

C. Szegedy et al. ICLR, 2014



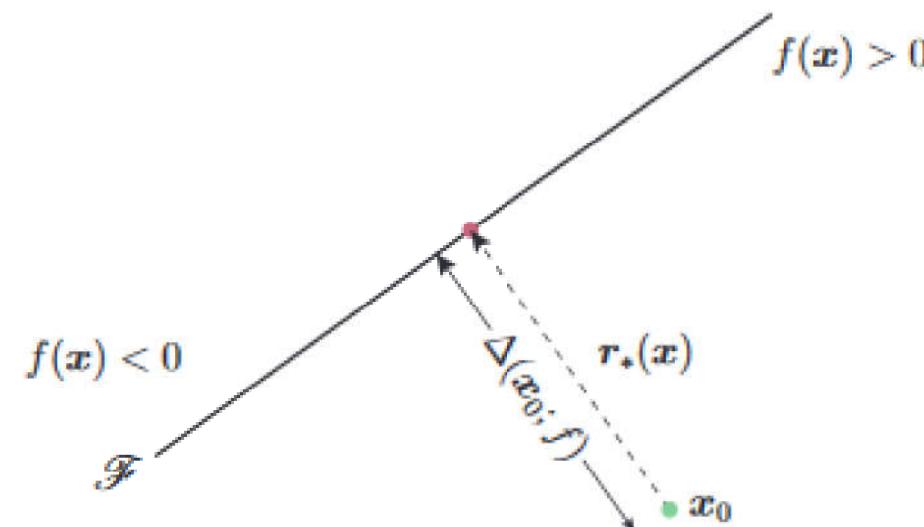
Explaining and Harnessing Adversarial Example

I.J. Goodfellow et al. ICLR, 2015



Deep Fool: a simple and accurate method to fool deep neural networks

S.M. Moosavi-dezfooli et al. CVPR, 2016



Algorithm 1 DeepFool for binary classifiers

```
1: input: Image  $x$ , classifier  $f$ .  
2: output: Perturbation  $\hat{r}$ .  
3: Initialize  $x_0 \leftarrow x$ ,  $i \leftarrow 0$ .  
4: while  $\text{sign}(f(x_i)) = \text{sign}(f(x_0))$  do  
5:    $r_i \leftarrow -\frac{f(x_i)}{\|\nabla f(x_i)\|_2^2} \nabla f(x_i)$ ,  
6:    $x_{i+1} \leftarrow x_i + r_i$ ,  
7:    $i \leftarrow i + 1$ .  
8: end while  
9: return  $\hat{r} = \sum_i r_i$ .
```



Adversarial timeline ...



Interiguing properties of Neural Network

C. Szegedy et al. ICLR, 2014



Explaining and Harnessing Adversarial Example

I.J. Goodfellow et al. ICLR, 2015



Deep Fool: a simple and accurate method to fool deep neural networks

S.M. Moosavi-dezfooli et al. CVPR, 2016

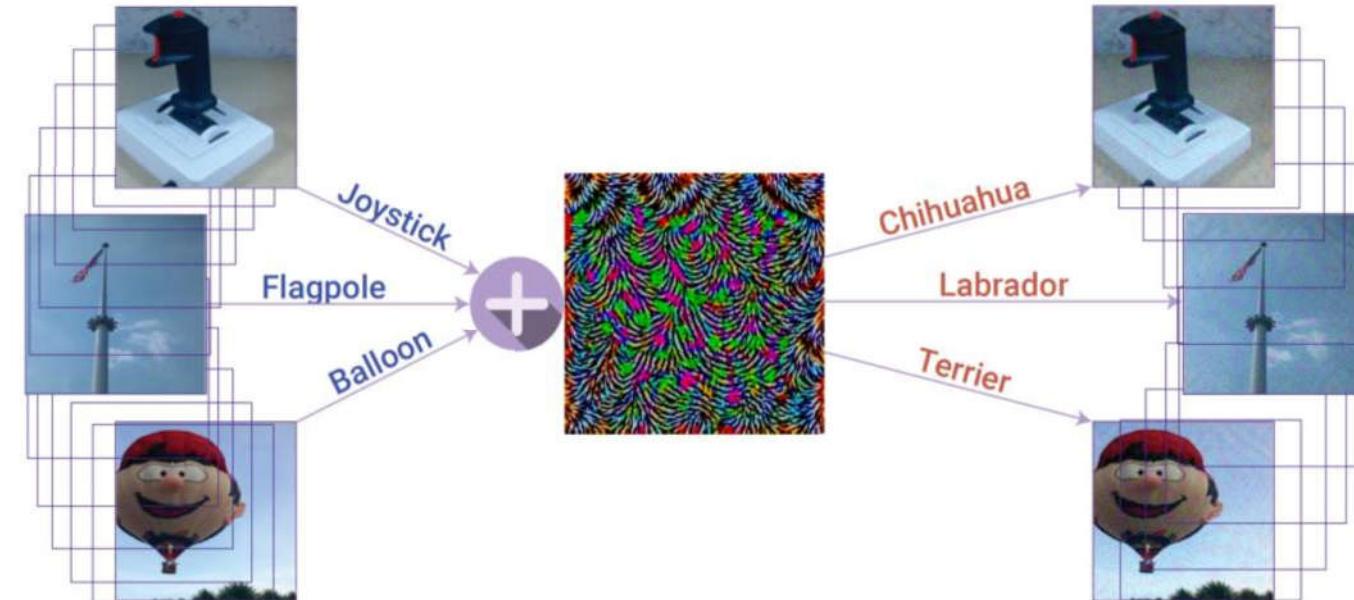


Universal adversarial perturbations

S.M. Moosavi-dezfooli et al. CVPR, 2017

$$\mathbb{P}_{x \sim \mu} \left(\hat{k}(x + v) \neq \hat{k}(x) \right) \geq 1 - \delta$$

$$\|v\|_p \leq \xi$$





Interiguing properties of Neural Network

C. Szegedy et al. ICLR, 2014



Explaining and Harnessing Adversarial Example

I.J. Goodfellow et al. ICLR, 2015



Deep Fool: a simple and accurate method to fool deep neural networks

S.M. Moosavi-dezfooli et al. CVPR, 2016



Universal adversarial perturbations

S.M. Moosavi-dezfooli et al. CVPR, 2017



Towards evaluating the robustness of neural networks

Carlini et al. Security and Privacy (S&P). 2017



Interiguing properties of Neural Network

C. Szegedy et al. ICLR, 2014



Explaining and Harnessing Adversarial Example

I.J. Goodfellow et al. ICLR, 2015



Deep Fool: a simple and accurate method to fool deep neural networks

S.M. Moosavi-dezfooli et al. CVPR, 2016



Universal adversarial perturbations

S.M. Moosavi-dezfooli et al. CVPR, 2017



Towards evaluating the robustness of neural networks

Carlini et al. Security and Privacy (S&P). 2017



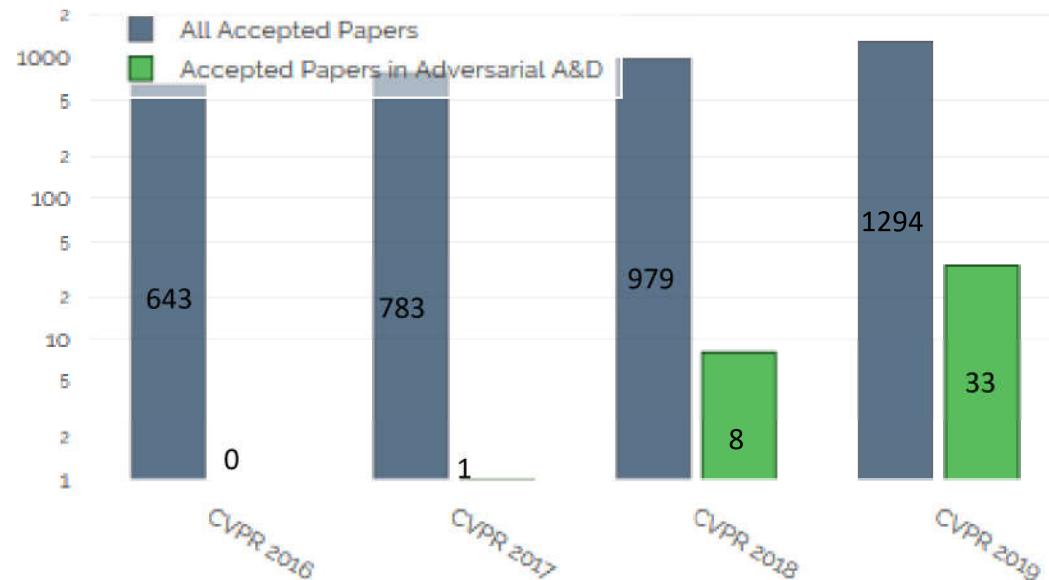
SparseFool: a few pixels make a big difference

Modas et al. CVPR, 2019

⋮

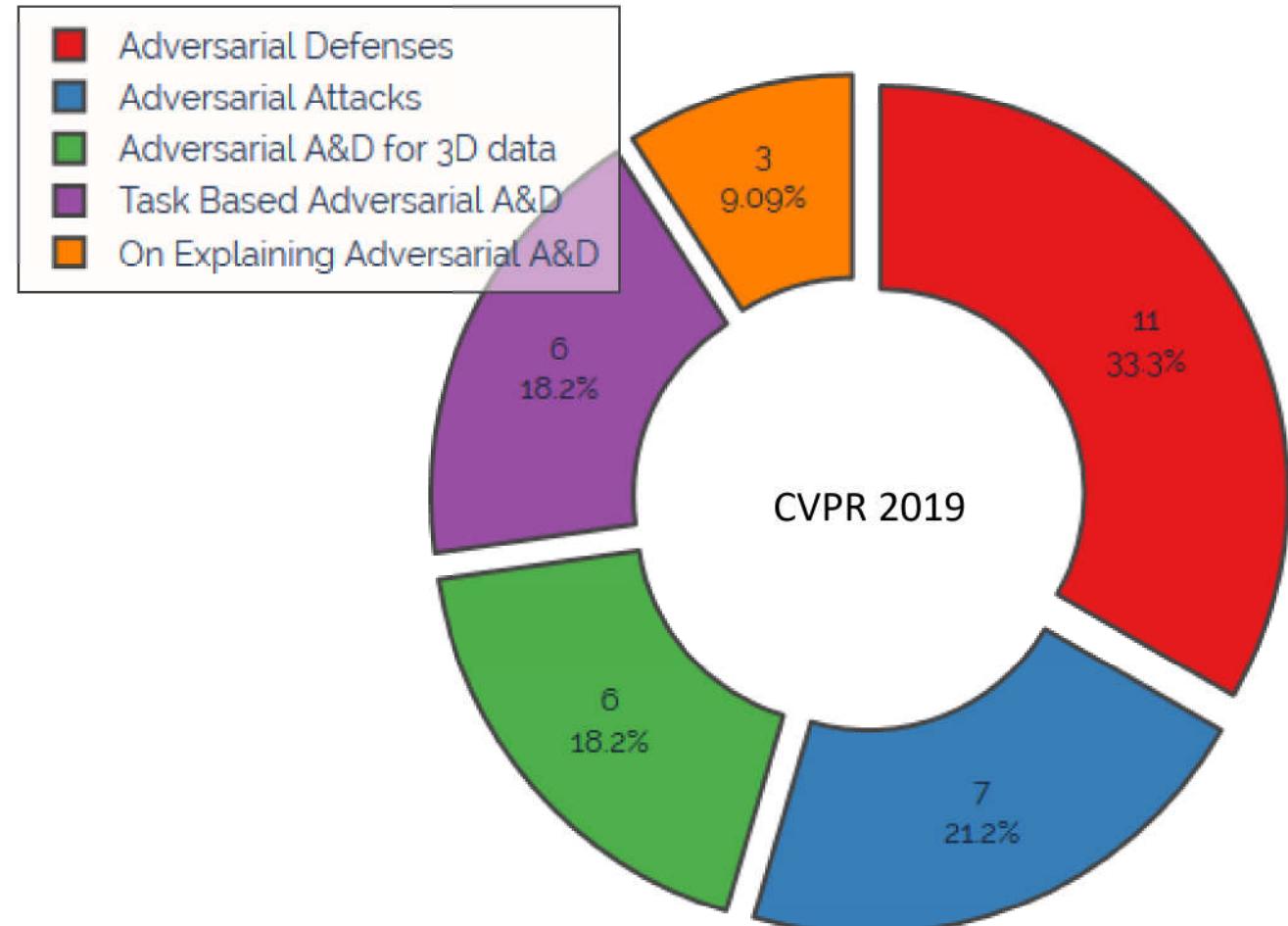
Increasing papers on Adversarial Attacks & Defenses

Accepted Papers at CVPR over the years



The number of accepted papers at CVPR over the years. Note the sharp rise in accepted papers in Adversarial A&D.

Distribution of Paper Topics



<http://fna.ir/a5d>

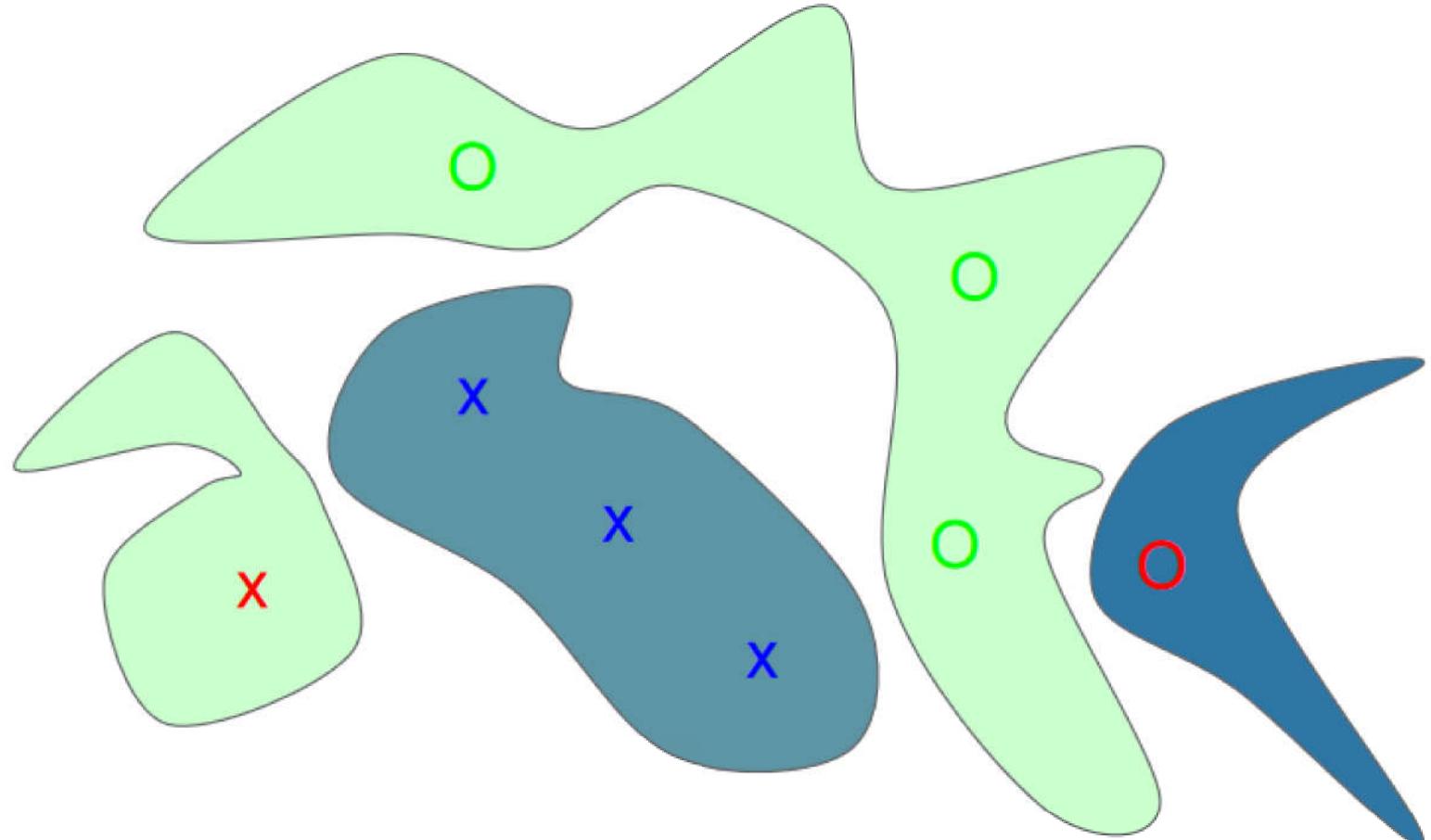
Various methods of attacks





Adversarial Examples from Overfitting

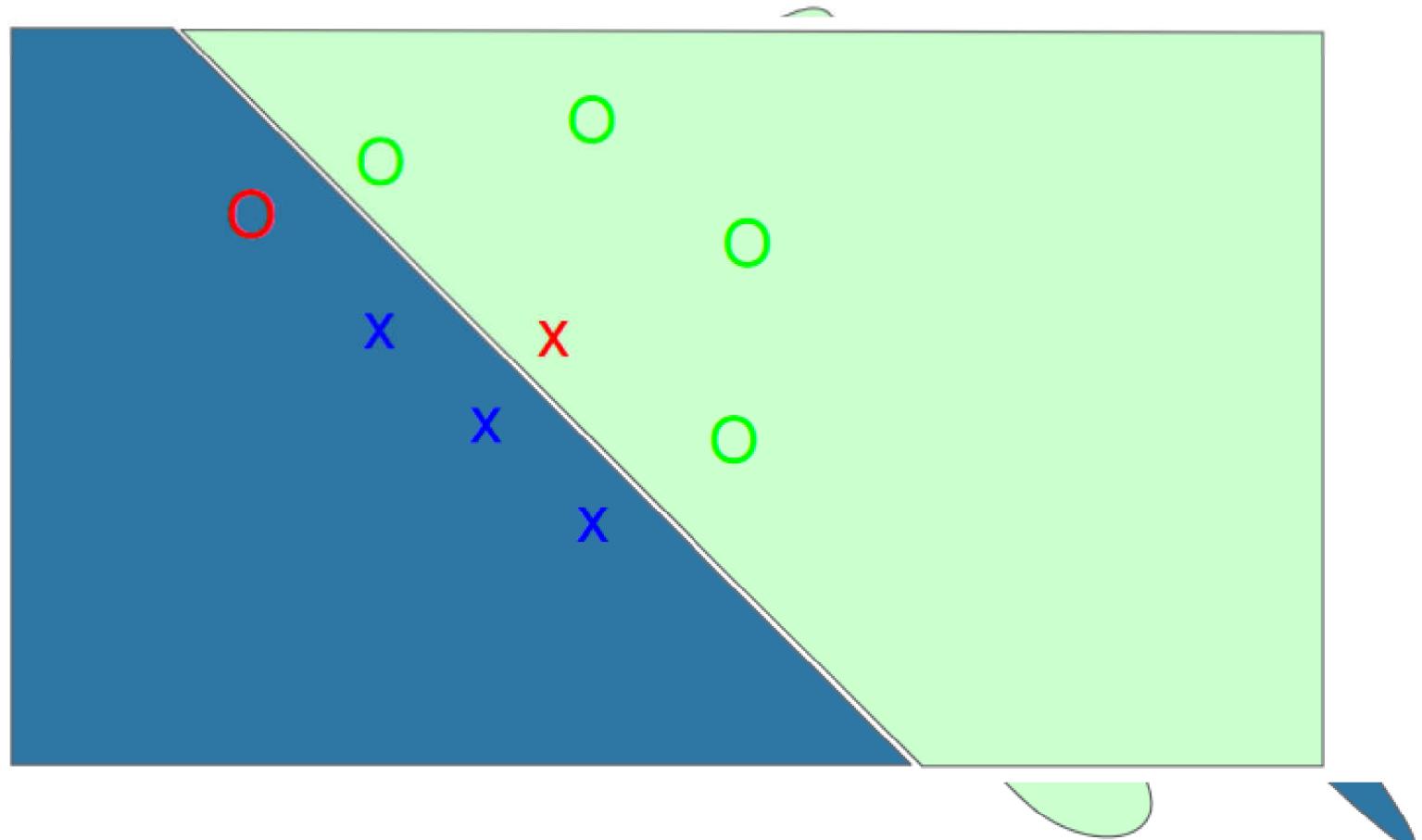
- Adversarial Examples rooted in :
 - **Overfitting**
 - Excessive Linearity



(Goodfellow 2016)



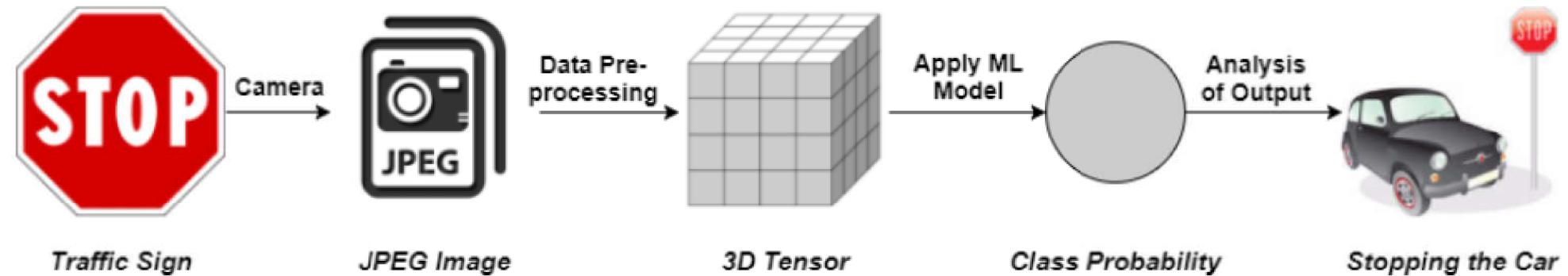
- Adversarial Examples rooted in :
 - Overfitting
 - Excessive Linearity



(Goodfellow 2016)

- The Attack Surface
- The Adversarial Capabilities
- The Adversarial Goals

- The Attack Surface



- Evasion Attack :: during the testing phase (* the most common type of attack!)
- Poisoning Attack :: during the training time
- Exploratory Attack :: during the testing phase (Given black box access to the model try to gain as much knowledge as possible)

- The Adversarial Capabilities

- Training Phase Capabilities

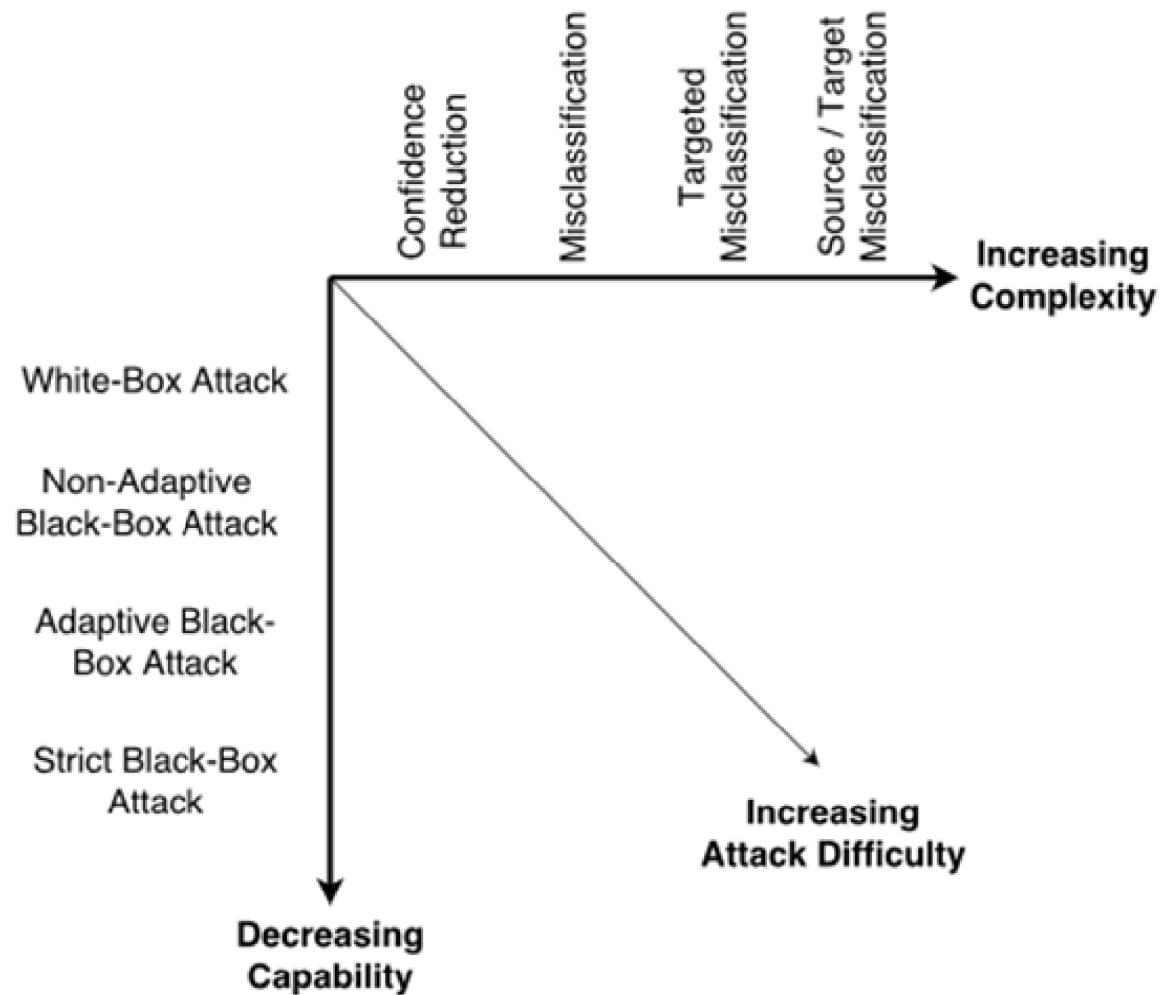
- Data Injection :: does not have any access to the training data as well as to the learning algorithm but has ability to augment a new data to the training set
- Data Modification :: does not have access to the learning algorithm but has full access to the training data
- Logic Corruption :: meddle with the learning algorithm

- Testing Phase Capabilities

- White-Box Attacks :: an adversary has total knowledge about the model (f), algorithm (train), training data distribution (μ), parameters (θ) of the fully trained model architecture
- Black-Box Attacks :: no knowledge about the model and uses information about the settings or past inputs
 - Non-Adaptive Black-Box Attack :: only gets access to the target model's training data distribution (μ)
 - Adaptive Black-Box Attack :: doesn't have any information regarding the training process but can access the target model as an oracle
 - Strict Black-Box Attack :: may not contain the data distribution(μ) but has the ability to collect the input-output pairs(x,y) from the target classifier. However, he can not change the inputs to observe the changes in output like an adaptive attack procedure

- **Adversarial Goals:**

- Confidence Reduction
 - The adversary tries to reduce the confidence of prediction for the target model
- Misclassification
 - The adversary tries to alter the output classification of an input example to **any** class different from the original class.
- Targeted Misclassification
 - The adversary tries to produce inputs that force the output of the classification model to be a **specific** target class
- Source/Target Misclassification
 - The adversary attempts to force the output of classification for a **specific** input to be a **particular** target class



Attack Difficulty with respect to adversarial capabilities and goals for Evasion Attacks





Adversarial example crafting procedures :

1. Direction sensitivity estimation

- The adversary evaluate the sensitivity of a class change to each input feature
- By identifying the direction in the data manifold around the example X

$$X_* = X + \arg \min_{\delta X} \{ \|\delta X\| : F(X + \delta X) \neq F(X) \}$$

Most DNN models make this formulation non-linear and non-convex, making it hard to find a closed-solution in most of the cases

2. Perturbation selection

- The adversary exploits the knowledge of sensitive information to select a perturbation
- Selecting the perturbation δX

3. Replace X with $X + \delta X$



Different techniques to find an approximate solution for adversarial example finding problem : **(Direction Sensitivity Estimation)**

- L-BFGS :: Szegedy et al.

$$\arg \min_r f(x + r) = l \quad \text{s.t. } (x + r) \in D$$

- Fast Gradient Sign Method (FGSM) :: Goodfellow et al

$$X_* = X + \epsilon * sign(\nabla_x J(X, y_{true}))$$

- Target Class Method :: Kurakin et al.

$$X_* = X - \epsilon * sign(\nabla_x J(X, y_{target}))$$

- Basic Iterative Method :: Kurakin et al.

$$X_*^0 = X; \quad X_*^{n+1} = Clip_{X,e}\{X_*^n + \alpha * sign(\nabla_x J(X_*^n, y_{true}))\}$$

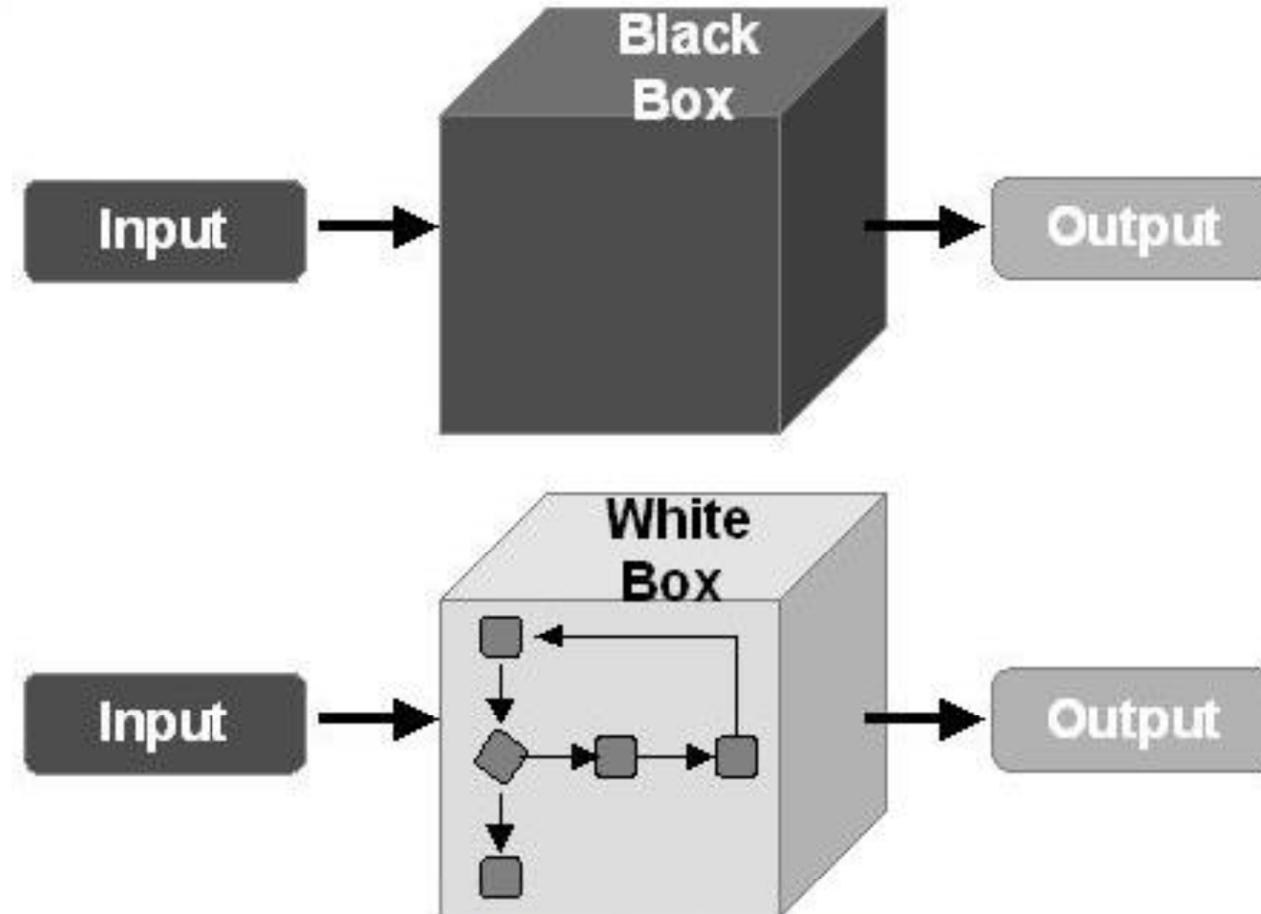
- Jacobian Based Method :: Papernot et al.

- Finding sensitivity direction by using **forward derivative**, which is the Jacobian of the trained model F

saliency value $S(x, t)[i] = \begin{cases} 0, & \text{if } \frac{\partial F_t}{\partial x_i}(x) < 0 \text{ or } \sum_{j \neq t} \frac{\partial F_j}{\partial x_i}(x) > 0 \\ \frac{\partial F_t}{\partial x_i}(x) \left| \sum_{j \neq t} \frac{\partial F_j}{\partial x_i}(x) \right|, & \text{otherwise} \end{cases}$

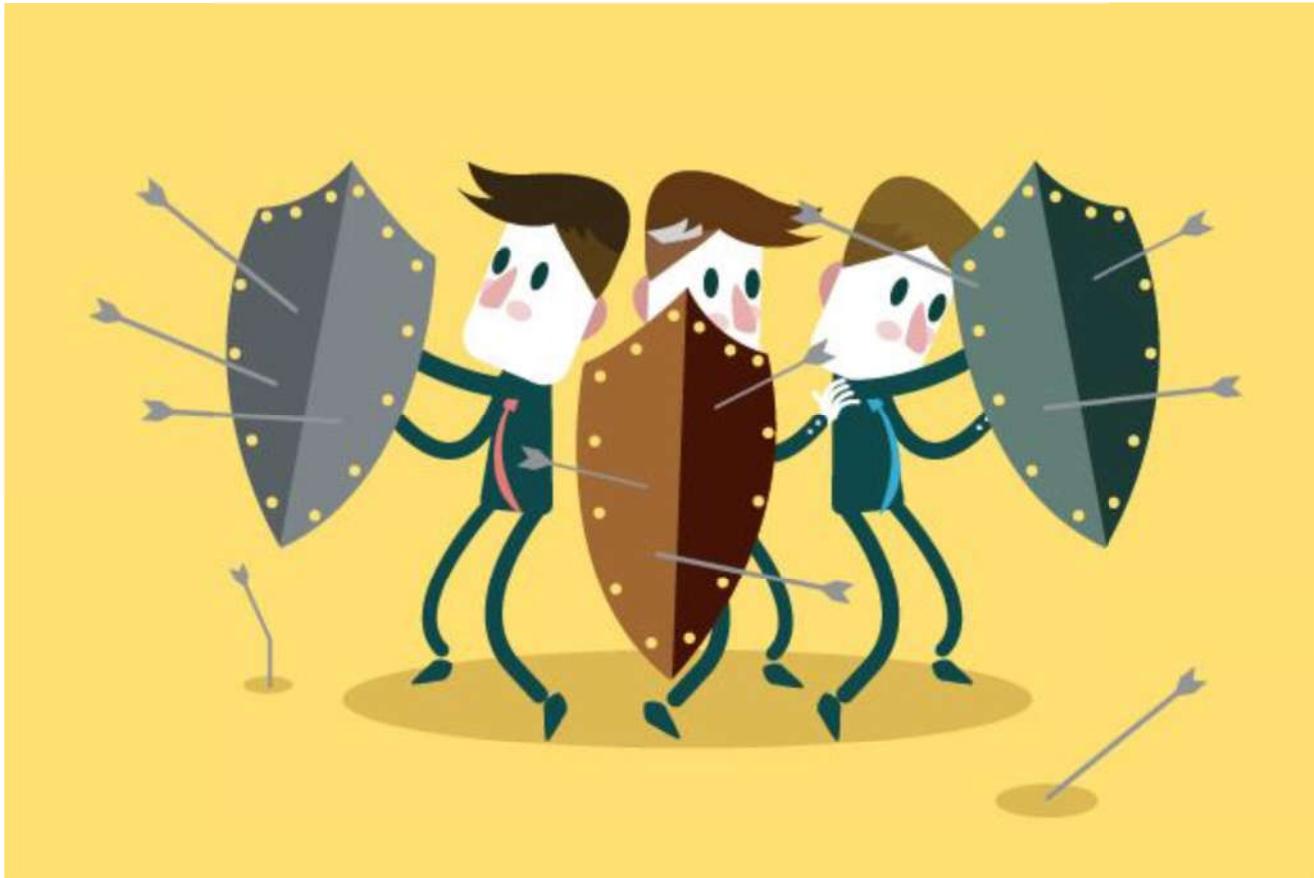


Crafting adversarial examples:





Defenses



<http://www.rmmagazine.com/2016/06/01/the-3-lines-of-defense-for-good-risk-management/>



□ Defense is hard!

*A theoretical model of the adversarial example crafting process is very **difficult** to construct.*

- Non-linearity
- non-convex
- Complex optimization process
- ...

- Most of the current defense strategies are
 - not adaptive to **all types** of adversarial attack
- Implementation of such defense strategies
 - may incur **performance overhead**



- Generative pretraining
- Removing perturbation with an autoencoder
- Adding noise at test time
- Ensembles
- Confidence-reducing perturbation at test time
- Dropout
- Adding noise at train time
- Various non-linear units
- ...





$\text{loss}(x, y)$



Small when prediction is
correct on legitimate input

$$\text{loss}(x, y) + \text{loss}(x + \epsilon \cdot \text{sign}(\text{grad}), y)$$



Small when prediction is
correct on legitimate input



Small when prediction is
correct on adversarial input

□ The existing defense mechanisms can be categorized among the following types based on their methods of implementation:

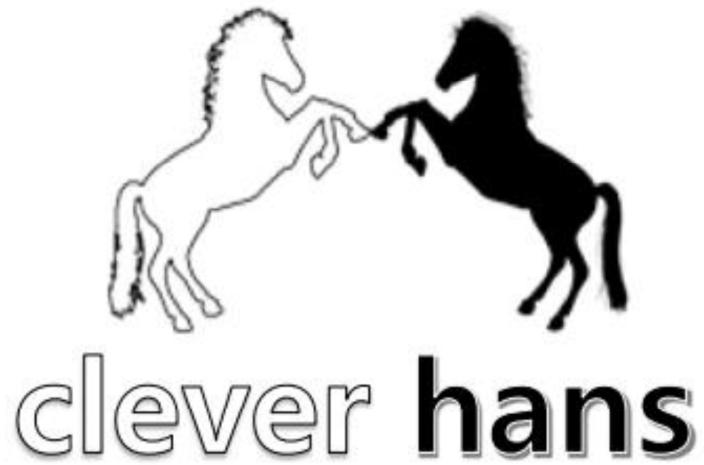
- Adversarial Training
- Gradient Hiding
- Defensive Distillation
- Feature Squeezing
- Blocking the Transferability
- Defense-GAN
- Mag Net
- Using High-Level Representation Guided De-noiser





ToolBox	Base Lib.	Usability	Updating
Clever-Hans	TensorFlow, Pytorch	Semi	Well
Fool-Box	TensorFlow, Keras, Pytorch, MXnet	Easy	Well
IBM ART	python	Semi	Well

...



<https://github.com/tensorflow/cleverhans>

<http://www.cleverhans.io/>

CleverHans Documentation

This documentation is auto-generated from the docstrings of modules of the current *master* branch of [tensorflow/cleverhans](#).

To get started, we recommend reading the [github readme](#). Afterwards, you can learn more by looking at the following modules:

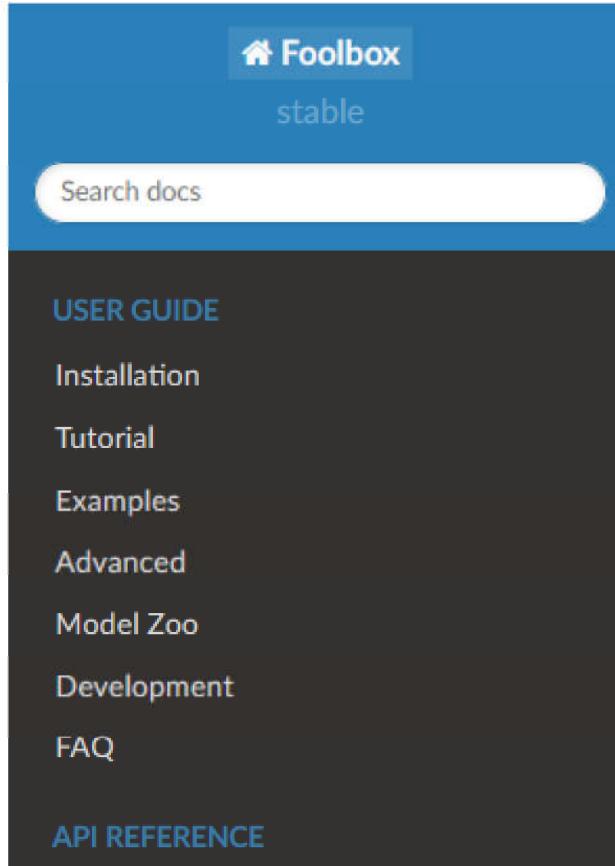
- [attacks module](#)
- [model module](#)

Indices and tables

- [Index](#)
- [Module Index](#)
- [Search Page](#)

<https://cleverhans.readthedocs.io/en/latest/>





The screenshot shows the homepage of the Foolbox documentation. The top navigation bar includes a 'Foolbox' icon, the text 'stable', and a search bar labeled 'Search docs'. Below this is a sidebar with 'USER GUIDE' and links to 'Installation', 'Tutorial', 'Examples', 'Advanced', 'Model Zoo', 'Development', and 'FAQ'. At the bottom of the sidebar is an 'API REFERENCE' section. The main content area has a breadcrumb trail 'Docs » Welcome to Foolbox' and a 'Edit on GitHub' button. The main title 'Welcome to Foolbox' is prominently displayed. A paragraph describes Foolbox as a Python toolbox for creating adversarial examples. It then states that it supports various frameworks: TensorFlow, PyTorch, Theano, Keras, Lasagne, and MXNet. A bulleted list on the right side details these frameworks.

Docs » Welcome to Foolbox [Edit on GitHub](#)

Welcome to Foolbox

Foolbox is a Python toolbox to create adversarial examples that fool neural networks.

It comes with support for many frameworks to build models including

- TensorFlow
- PyTorch
- Theano
- Keras
- Lasagne
- MXNet



Models

API Reference

- [foolbox.models](#)
 - [Models](#)
 - [Wrappers](#)
 - [Detailed description](#)
- [foolbox.criteria](#)
 - [Criteria](#)
 - [Examples](#)
 - [Detailed description](#)
- [foolbox.zoo](#)
 - [Get Model](#)
 - [Fetch Weights](#)
- [foolbox.distances](#)
 - [Distances](#)
 - [Aliases](#)
 - [Base class](#)
 - [Detailed description](#)
- [foolbox.attacks](#)
- [foolbox.adversarial](#)
- [foolbox.utils](#)

Model	Base class to provide attacks with a unified interface to models.
DifferentiableModel	Base class for differentiable models that provide gradients.
TensorFlowModel	Creates a Model instance from existing TensorFlow tensors.
TensorFlowEagerModel	Creates a Model instance from a TensorFlow model using eager execution.
PyTorchModel	Creates a Model instance from a PyTorch module.
KerasModel	Creates a Model instance from a Keras model.

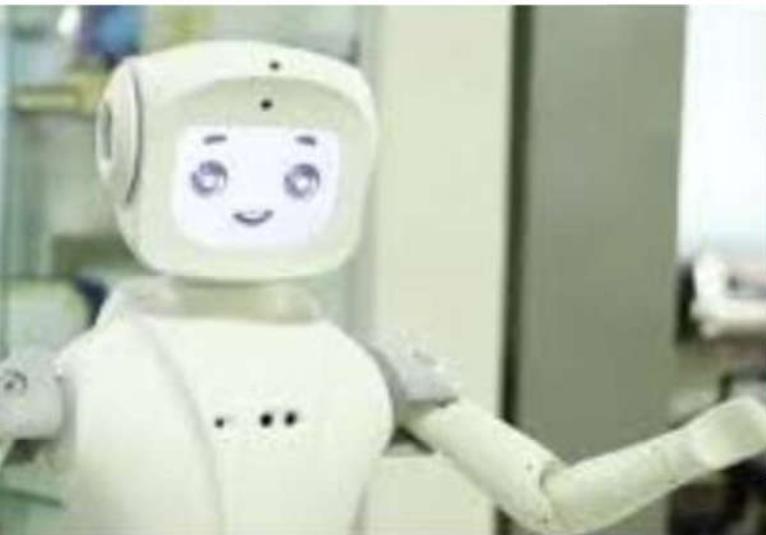
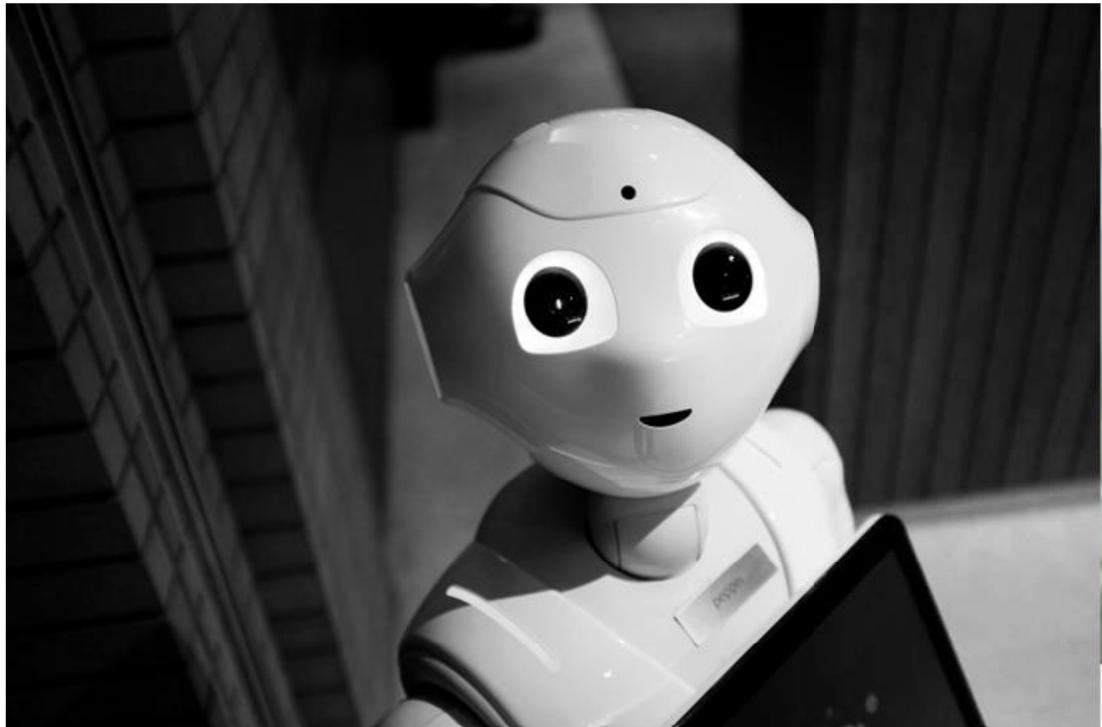
Criteria

We provide criteria for untargeted and targeted adversarial attacks.

Misclassification	Defines adversarials as images for which the predicted class is not the original one.
TopKMisclassification	Defines adversarials as images for which the original class is not one of the top k predicted classes.
OriginalClassProbability	Defines adversarials as images for which the probability of the original class is below a threshold.
ConfidentMisclassification	Defines adversarials as images for which the probability of any class other than the original one is above a threshold.
TargetClass	Defines adversarials as images for which the predicted class is the given target class.
TargetClassProbability	Defines adversarials as images for which the probability of a given target class is above a threshold.



Tips to stay safe





- Attacking is easy
- Defending is difficult
- Adversarial training provides regularization and semi-supervised learning
- The out-of-domain input problem is a bottleneck for model-based optimization generally
- There exist certain countermeasure, but none of them can act as a panacea for all challenges!



Mohammad Khalooei

Mkhalooei [at] gmail.com

Khalooei [at] aut.ac.ir

<https://ceit.aut.ac.ir/~khalooei>