



UNIVERSITÀ
DEGLI STUDI DI BARI
ALDO MORO



DIPARTIMENTO
DI INFORMATICA



Adversarial Machine Learning : an introduction

Dr. Giuseppina Andresini

Artificial Intelligence



- **AI is the new electricity**

- AI is the new electricity. It will transform every industry and create huge economic value.
- Technology like supervised learning is automation on steroids. It is very good at automating tasks and will have an impact on every sector – from healthcare to manufacturing, logistics and retail (cit. Andrew NG)



Deep Learning

Modern AI is:

Numerical Optimization + Data



Deep Learning, Machine Learning

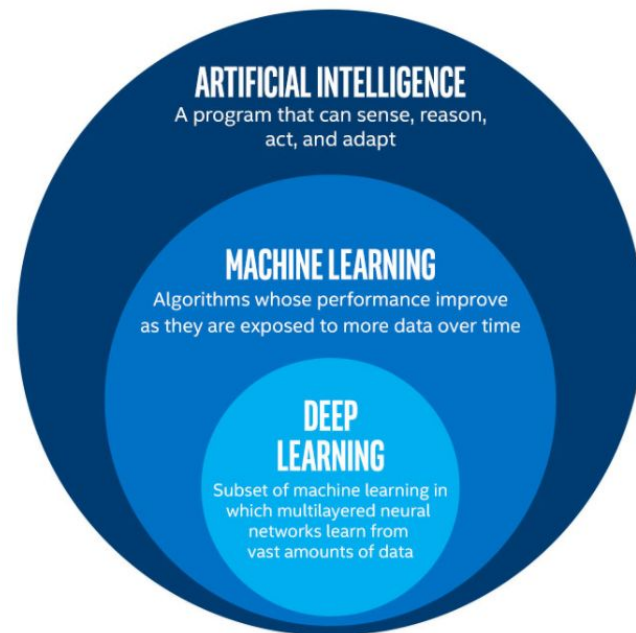
Deep Learning (DL) is:

1. A branch of **Machine Learning (ML)**
2. A set of algorithms that attempt to model **high-level abstractions** representation of data **from raw data**
3. Inspired by the function of the **biological brain**
4. Composed by multiple processing layers of multiple non-linear transformations

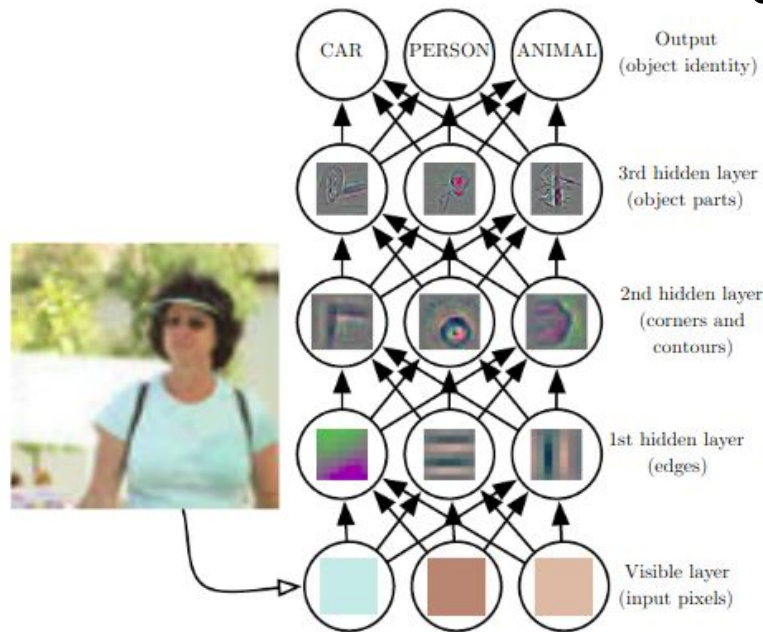
DL is a branch of ML



- Data preparation:
 - ML requires **structured data as input**
 - Human intervention to identify the features based on the data type
 - e.g., for images pixel value, shape, rotation
 - DL extract the high-level abstractions representation of data **from raw data** without human intervention
- Approach
 - ML need to human intervention to provide adjustments to improve the accuracy of the model
 - DL performs an **end-to-end learning** :
 - From raw data to a task (e.g., classification) the model learn how to do it automatically



End-to-end learning



- DL models extract learns representations (features) from the image in automatically manner.

- For each hidden layer, given the pixels:
 - the **1st hidden layer** identify edges, by comparing the brightness of neighboring pixels.
 - the **2nd hidden layer** search for corners and extended contours (as collections of edges)
 - the **3rd hidden layer** detect parts of specific objects (finding specific collections of contours and corners).
 - Finally, this description of the image in terms of the object parts it contains is used to **recognize the objects present in the image.**

Deep learning : a recap

BRACE YOURSELF

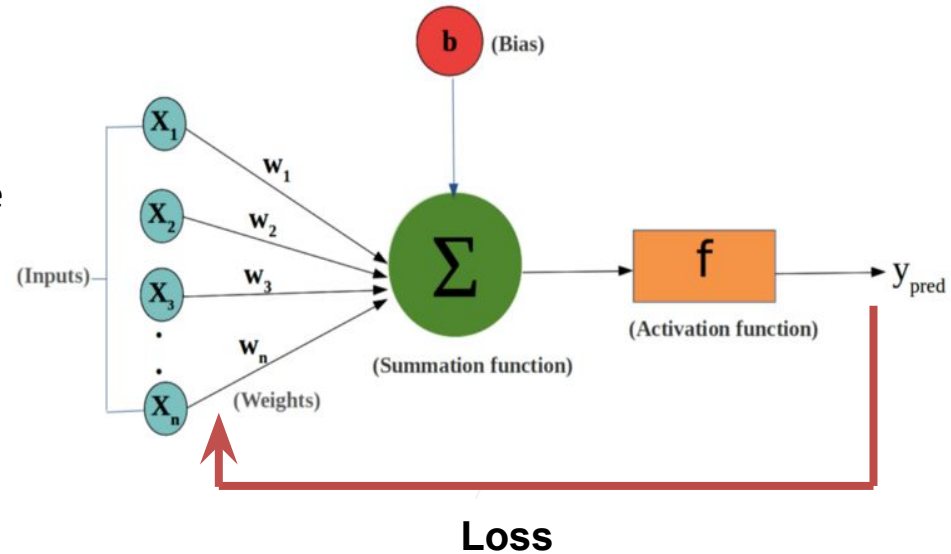


RECAP IS COMING!

How NNs learn?

Learning by **trial-and-error**

- **Trial:**
 - Compute the prediction of a given input data
- **Evaluate:**
 - Provide a feedback of the correctness of the prediction
 - Comparing the actual output with the expected input
 - $loss = y - \hat{y}$
- **Adjust:**
 - The weights are iteratively adjusted in order to reduce the error



The NN is refined over time to find the weights that minimise the error on the training and provide more accurate predictions

Loss function

- The loss function measures how much the predicted output label is different from the **expected output (ground truth label)**.
- The error at each hidden layer is an average of the evaluated error
- The choice of loss function is sensitive to the application task:

Mean squared error (MSE)

- The most commonly used loss function for regression
- Squared differences between true and predicted values

$$L_{\text{mse}}(y, \hat{y}) = ||y - \hat{y}||^2$$

Cross entropy loss

- Used in classification task
- Measures the performance of a classification model whose output is a probability between 0 and 1 by increasing its value when the predicted probability diverges from the actual label.

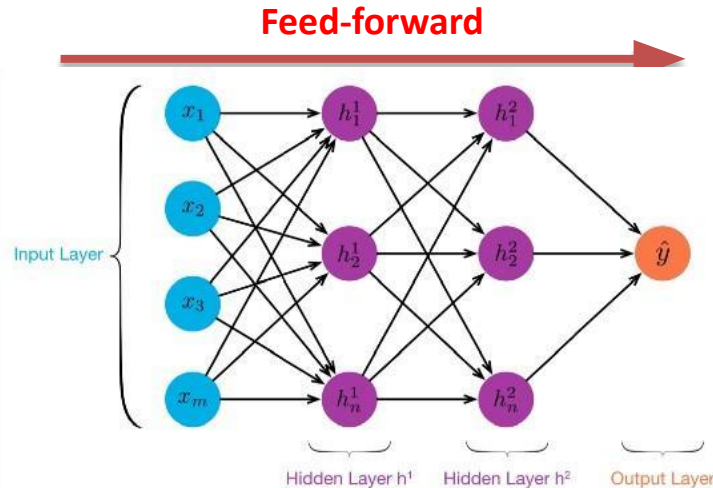
$$L_{\text{multiclasscrossentropy}}(y, p(\hat{y})) = - \sum_c y_c \log p(\hat{y})_c$$

Training a NN with backpropagation

- In **single-layer NNs** the error (loss) is computed as a direct function of weights
- In **multi-layer NNs** the loss is a composition function of the weights in earlier layers:
 - The gradient of composition function is performed by using the **backpropagation algorithm**:
 - Two phases: *forward* and *backward*

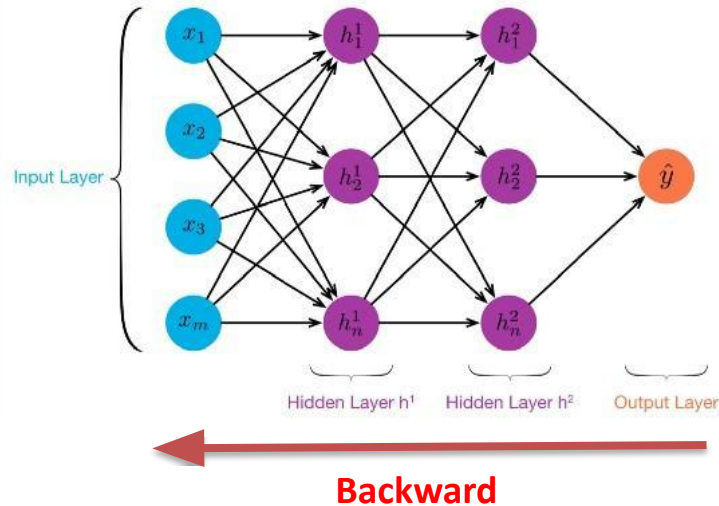
Training: forward phase

- The inputs are fed into NNs
- Computation across layers using the current set of weights
- Final prediction:
 - **Compare** the output with training
 - **Compute** the derivative of the loss with respect to the output



Training: backward phase

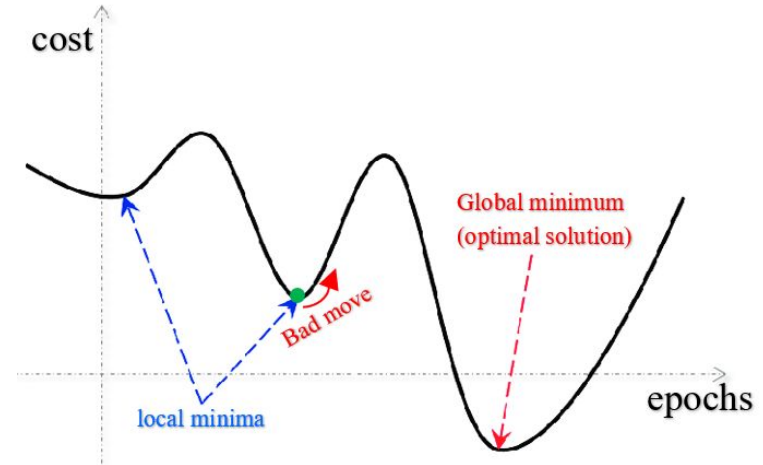
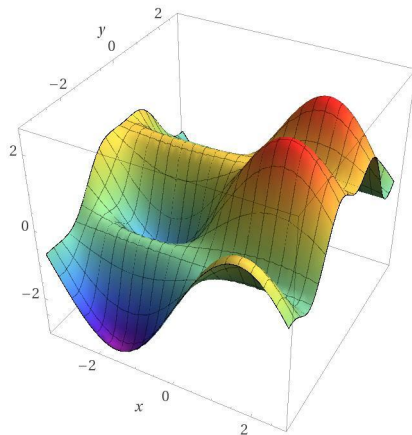
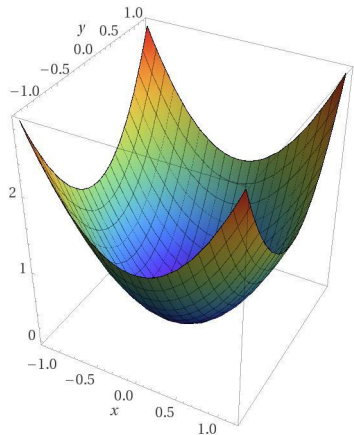
- Learn the **gradient** of a loss function with respect to the weights
 - Using the **chain rule**
 - The gradient is used to update the weights
 - Starting from output layer to input



Gradient Descent



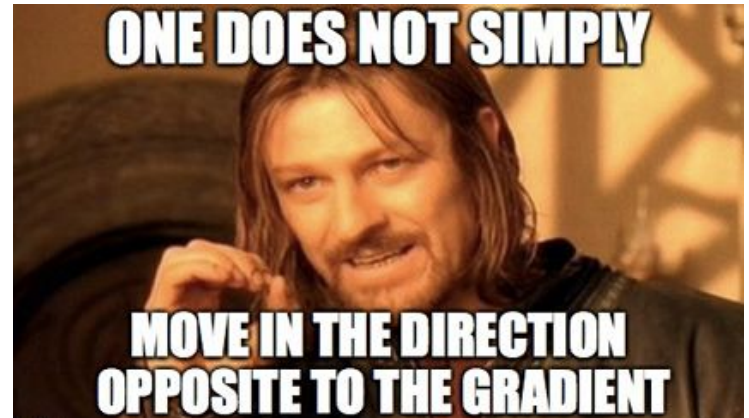
- First-order iterative optimization algorithm
 - Goal: find a **local minimum** of a differentiable function (loss function)
 - To train NNs and update the weights
 - To find the weights that minimize the cost function (loss function)
 - Partial derivative with y respect to its inputs x (dy/dx)



Gradient Descent

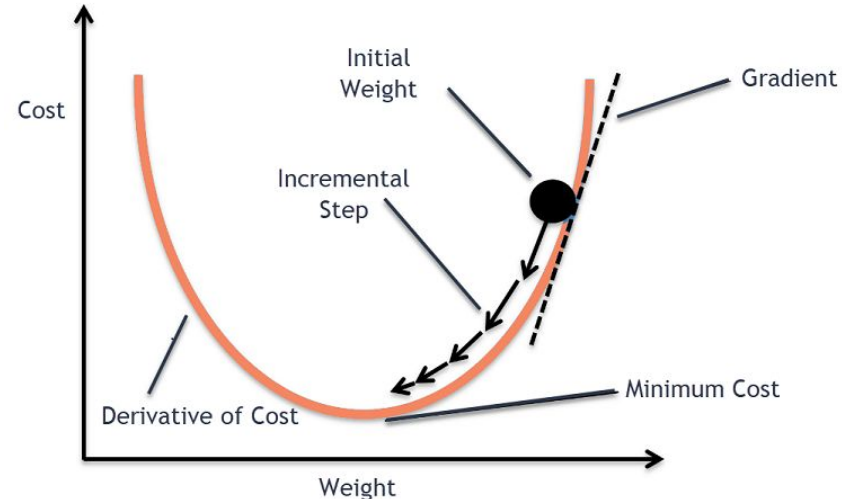


- A gradient measures how much the output of a function changes if you change the input
 - Measures the change in all weights with regard to the change in error.
 - Slope of a function
 - The steeper the slope and the faster a model can learn.
 - If the slope is zero, the model stops learning.



Gradient Descent

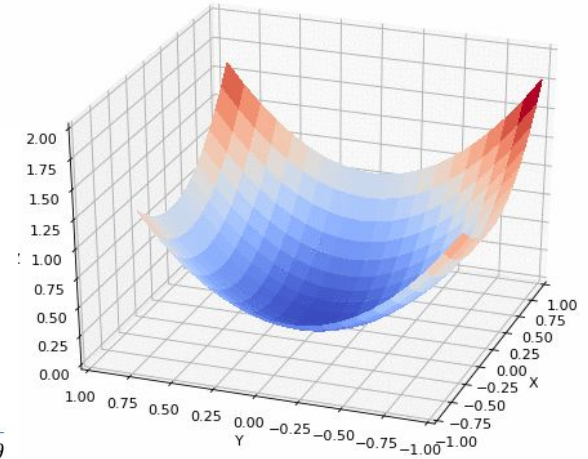
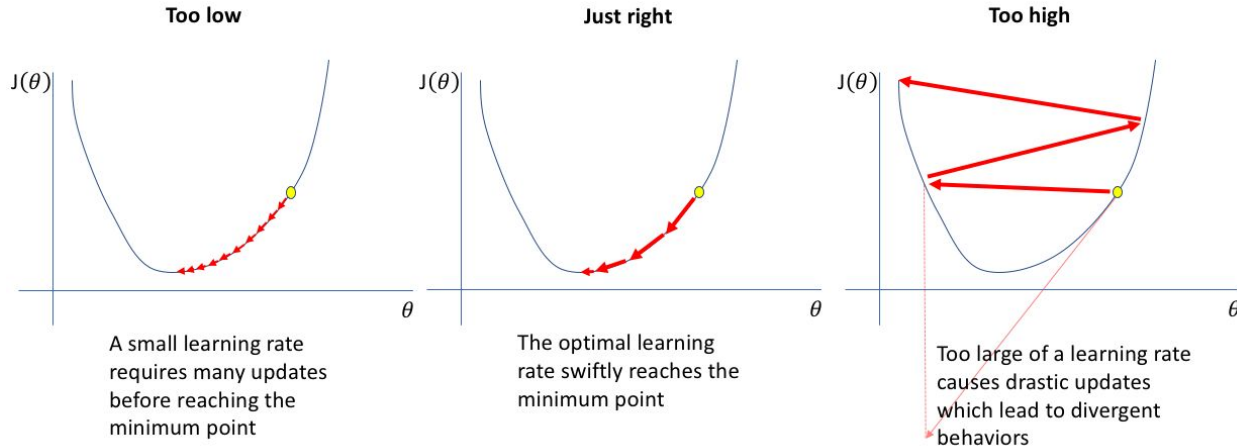
- Let us consider a **cost function**
 $y=f(x)=L(x,y)$
- **Goal: minimizing $f(x)$**
- On each iteration:
 - Compute the **slope (gradient)** of $f(x)$ at current point
 - Update the parameters in **the opposite direction** of the gradient
 - the parameters where the gradient gives the direction of the steepest ascent
 - by moving x in small steps with opposite sign of the derivative
 - the size of the step on each iteration is determined by the **learning rate α**
 - Until we reach a local minimum



$$W' = W - \alpha \nabla_W L(x, y, W)$$

Gradient Descent

- Learning rate: for gradient descent
 - scale the magnitude of our weight updates in order to minimize the network's loss function

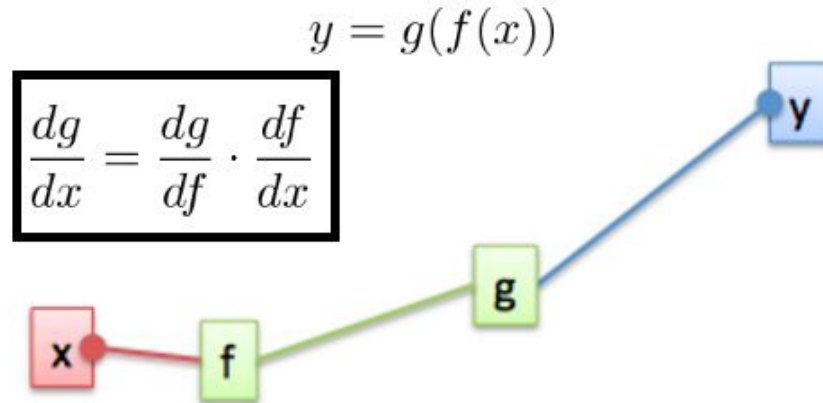


15

$$W' = W - \alpha \nabla_W L(x, y, W)$$

Chain rule

- In multi-layer NN the loss is a composition function of the weights in earlier layers
- The gradient is computed using backpropagation algorithm that leverages a **chain rules**
 - It calculates the error gradient in terms of summations of local-gradient products



source:
<https://bit.ly/3Gm5Ap5>

Object classification and object detection

Classification



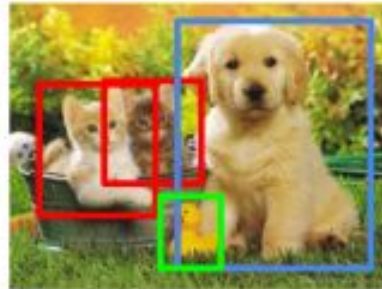
CAT

Classification
+ Localization



CAT

Object Detection



CAT, DOG, DUCK

Instance
Segmentation



CAT, DOG, DUCK

Single object

Multiple objects

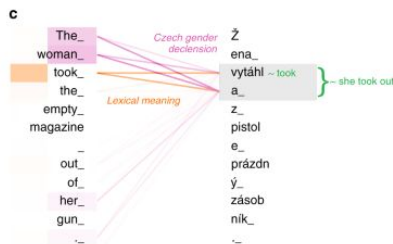
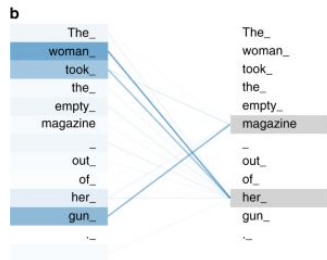
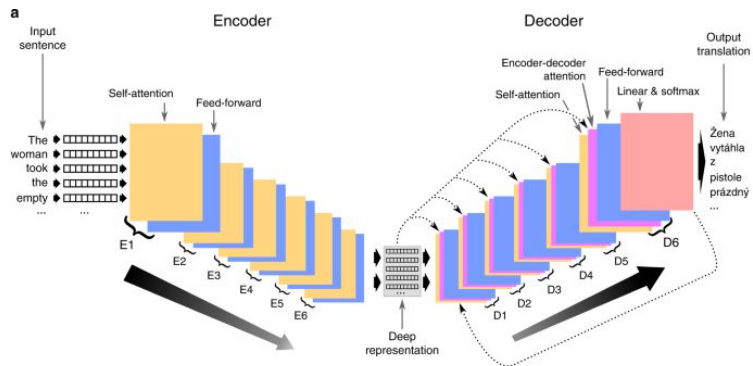
Is this a cat?

What is there in image and where?

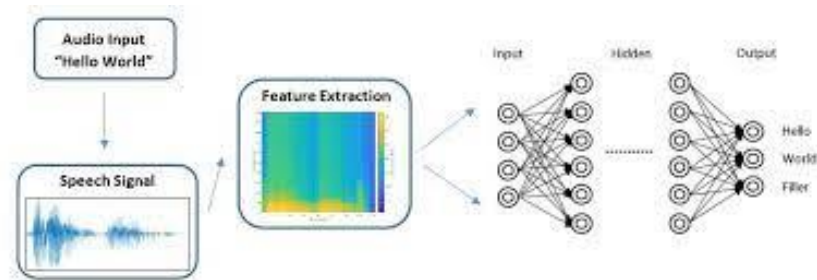
Which pixels belong to which object?

Text

Language Translation



Object recognition: speech

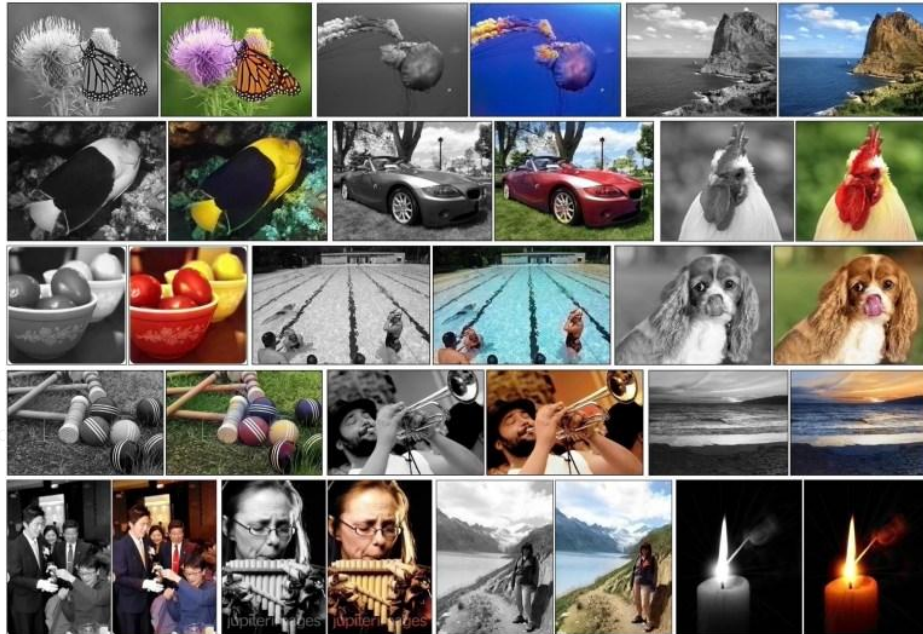


Automatic handwriting generation

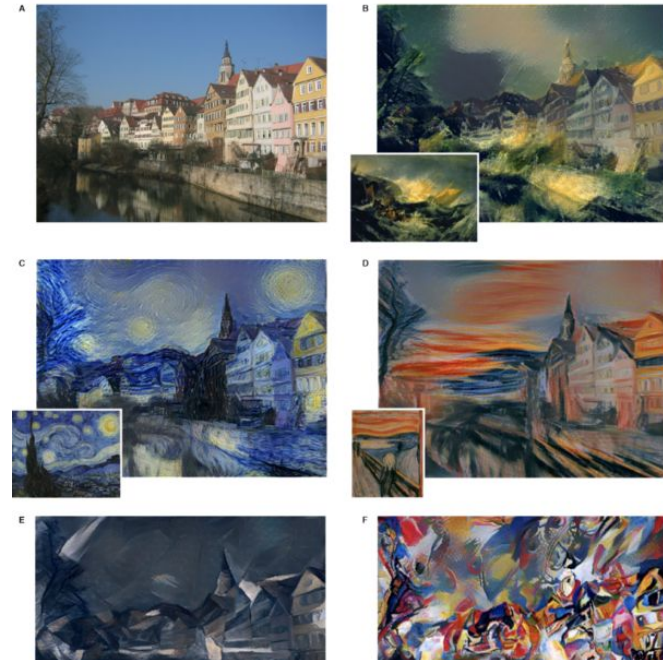
it was not written by me
it was not written by me
it was not written by me

Art

Colorize Black & White Images



Style transfer



Images

Faking faces and video

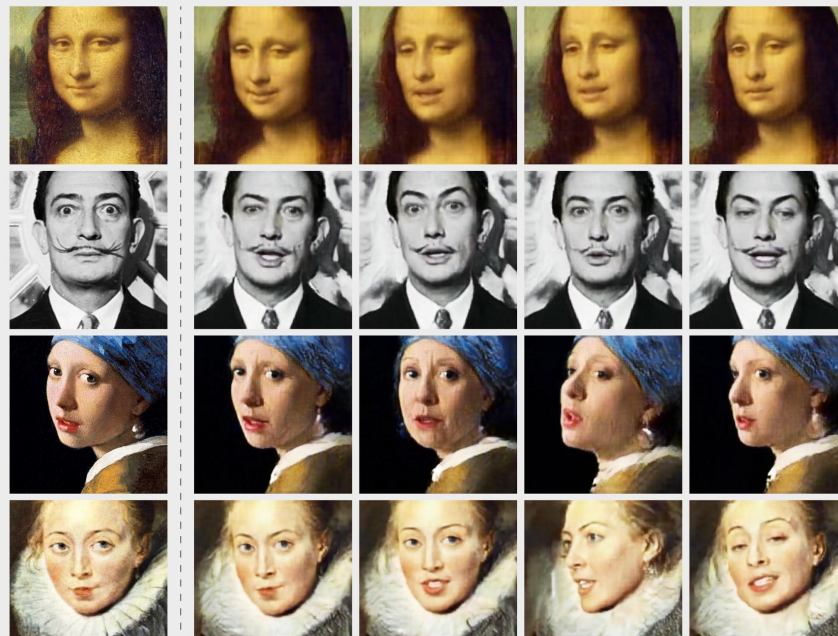


Original Video for Input Audio



b) Our result

Synthesizing Obama: Learning Lip Sync from Audio
<https://www.youtube.com/watch?v=9Yq67CjDqww>



Images

Image reconstruction



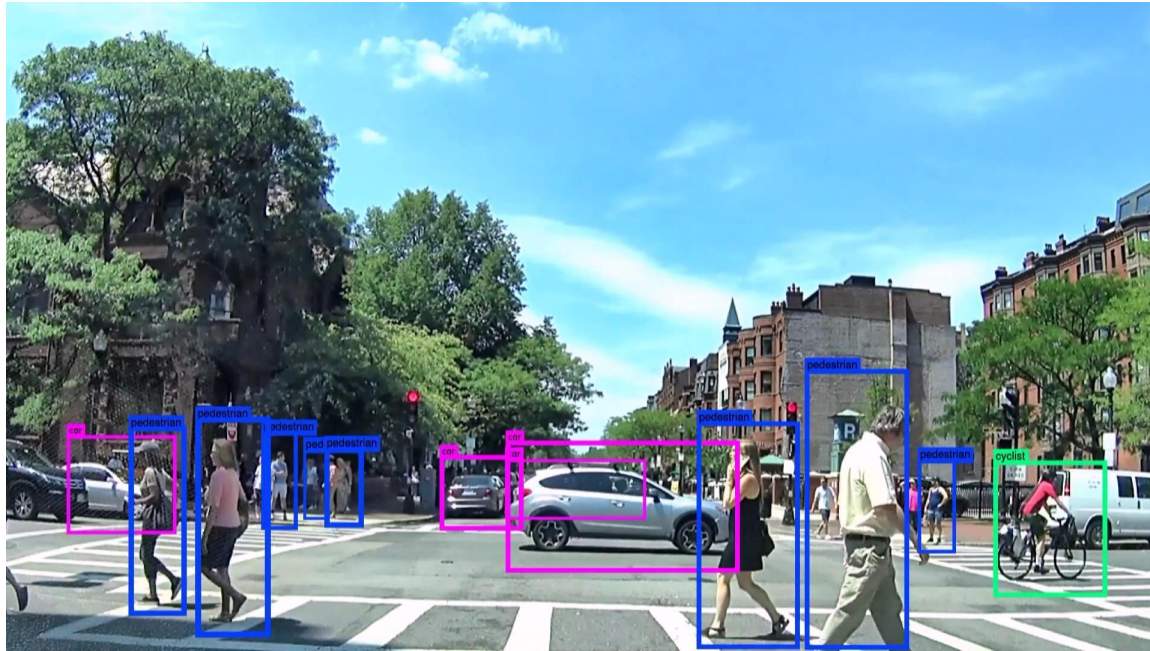
Images

Face creation



Self driving

Self driving cars



Cybersecurity



Cybersecurity: intrusion detection, malware detection, fraud detection...



Dark side of AI: adversarial learning



AI in Cybersecurity

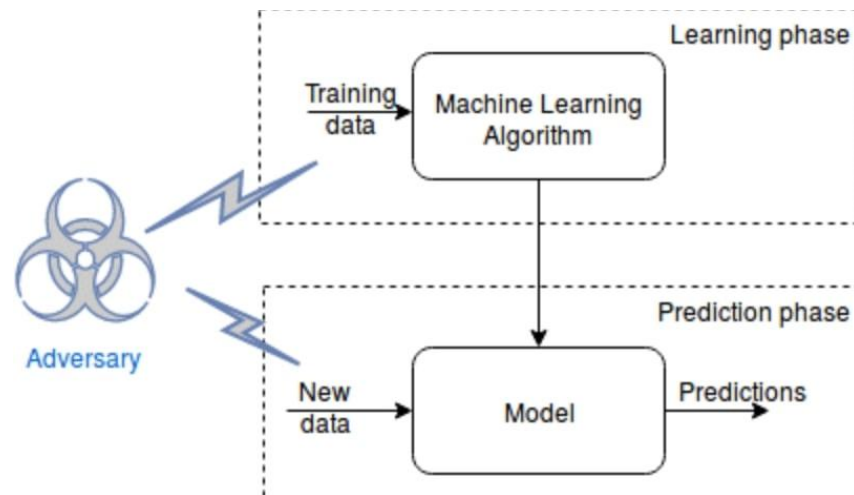


- Cybersecurity products are increasingly incorporating **Artificial Intelligence (AI)** technologies (e.g., Machine Learning, Deep Learning)
 - **PRO:**
 - Increase the **ability** to detect attacks
 - Reduce the **time** spent for threat detection and incident responses
 - **CONS:**
 - Attackers use AI to better understand their targets and design **new attacks**
 - AI systems are vulnerable to **adversarial attacks**



Adversarial machine learning

- Study of **attacks** on machine learning (and deep learning) algorithms
 - Collection of techniques to train neural networks on how to spot intentionally misleading data or behaviors
 - **Defenses** against such attacks
 - preemptively locate vulnerabilities and craft more flexible learning algorithms



Adversarial biometric

Chaos Computer Club (CCC), demonstrated that a fingerprint of the phone user, photographed from a glass surface, was enough to create a fake finger that could unlock an iPhone 5s secured with TouchID.

iPhone 5s fingerprint sensor 'hacked' within days of launch

A group of German hackers has found a way to bypass Apple's TouchID, and claims that fingerprint biometrics is an unsuitable method of access control.



Police are forced to destroy the fingerprints and DNA samples of terror suspects who are not charged Photo: ALAMY



By Sophie Curtis
9:33AM BST 23 Sep 2013

Follow 11.7K followers

The group, known as the Chaos Computer Club (CCC), demonstrated that a fingerprint of the phone user, photographed from a glass surface, was enough to create a fake finger that could unlock an iPhone 5s secured with TouchID.



Fingerprint images coming from a live (left side) and a fake finger (right side), almost indistinguishable.

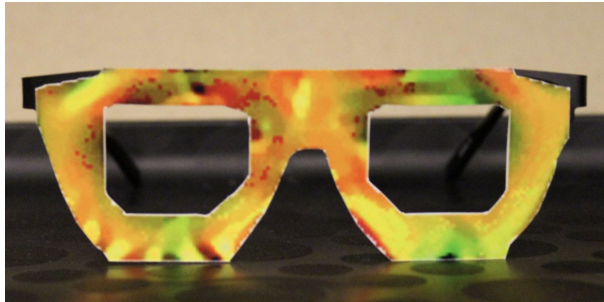


Fake fingers provided by PRA Lab for the [ICB 2013](#) – Spoofing Challenge

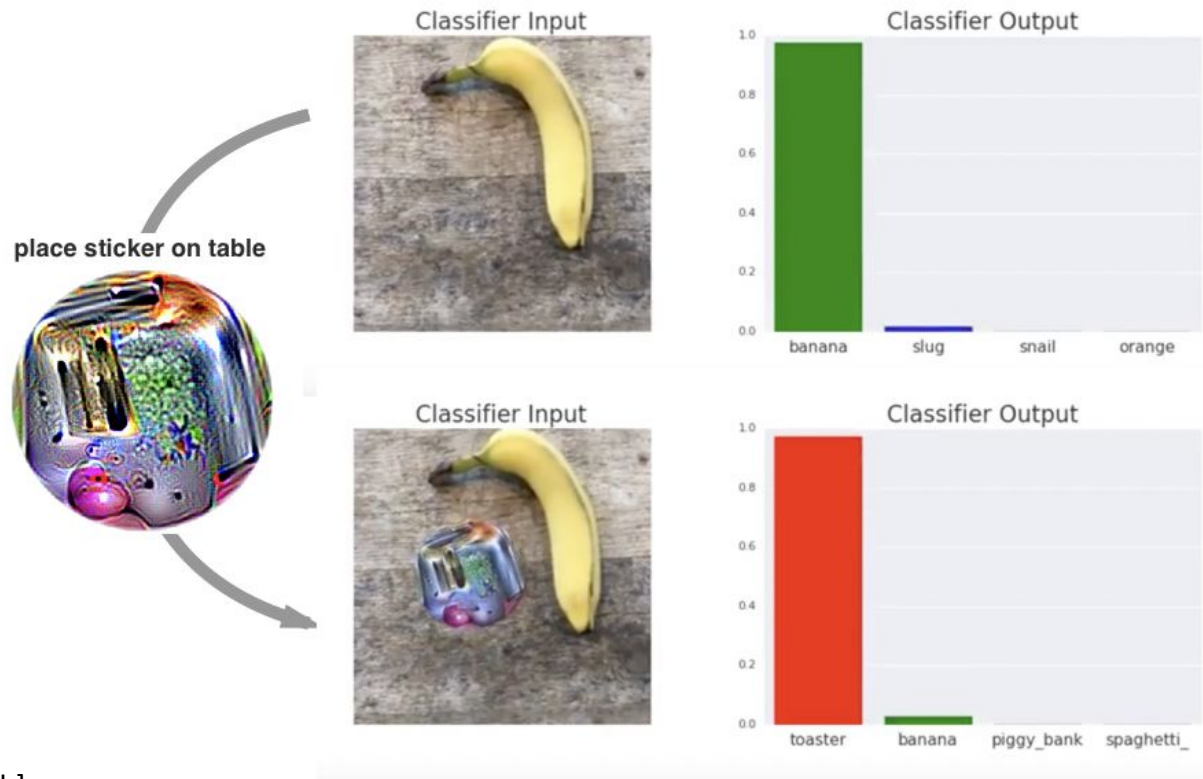


<https://www.euronews.com/2014/03/17/behind-the-mask-of-biometric-security?jwsource=cl>

Adversarial Image



Adversarial Image



[Brown et al 2017 -Adversarial Patch]

Adversarial audio



https://nicholas.carlini.com/code/audio_adversarial_examples/

Mozilla implementation of DeepSpeech.

[Reveal Transcription] “okay google browse to evil dot com”

[Reveal Transcription] “without the dataset the article is useless”

What are adversarial examples?



Puppy or Muffin?

SheepDog or Mop?

Adversarial examples are inputs to machine learning models that an attacker has intentionally designed to cause the model to make a mistake; they're like optical illusions for machines.

<https://openai.com/blog/adversarial-example-research/>

Categorization of adversarial attacks

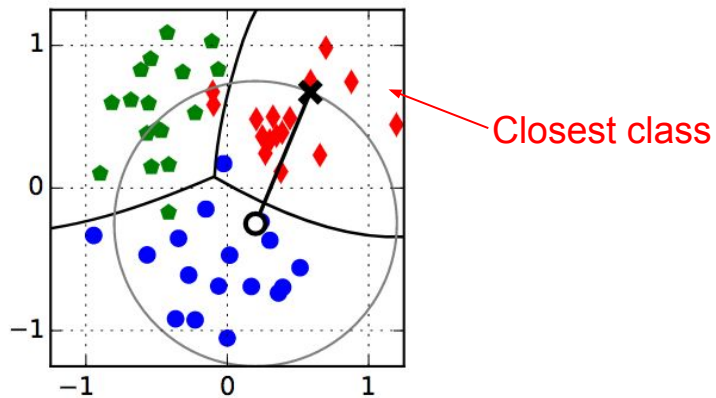


- **Goal:**
 - Error generic -> untarget attacks
 - Error specific -> target attacks
- **Capability**
 - Poison
 - Evasion
- **Knowledge**
 - White box
 - Black box

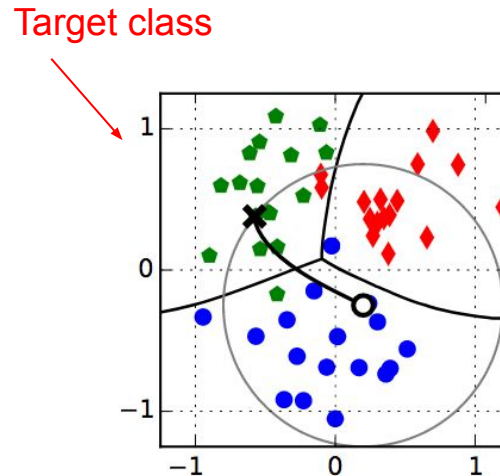
Attacker's goal

- **Error Specificity:**

- sample misclassified as a specific class or misclassified as any of the classes different from the true class



Error-generic

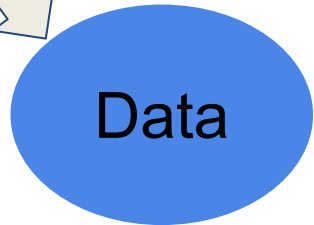


Error-specific

Attacker's capability

- **Poisoning:**
 - Training-time
- **Evasion:**
 - Test-time

Poisoning

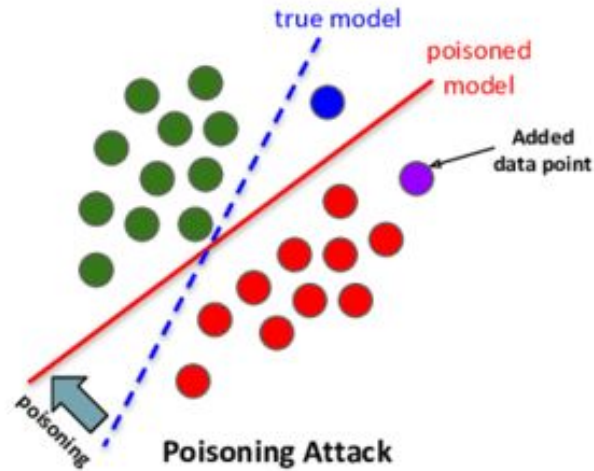


Evasion



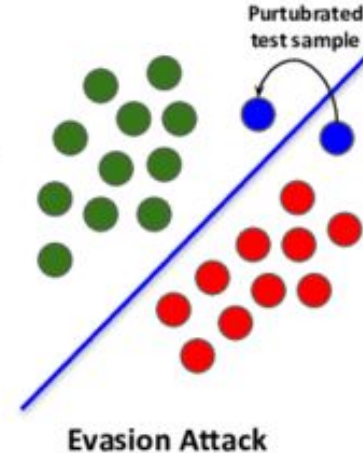
Attacker's capability

- Injection of bad data into the model's training



- A perturbed example is create to be misclassified by the model

- e.g., spam emails -try a number of different email contents against the model and discover a way to get their spam email classified as innocuous.

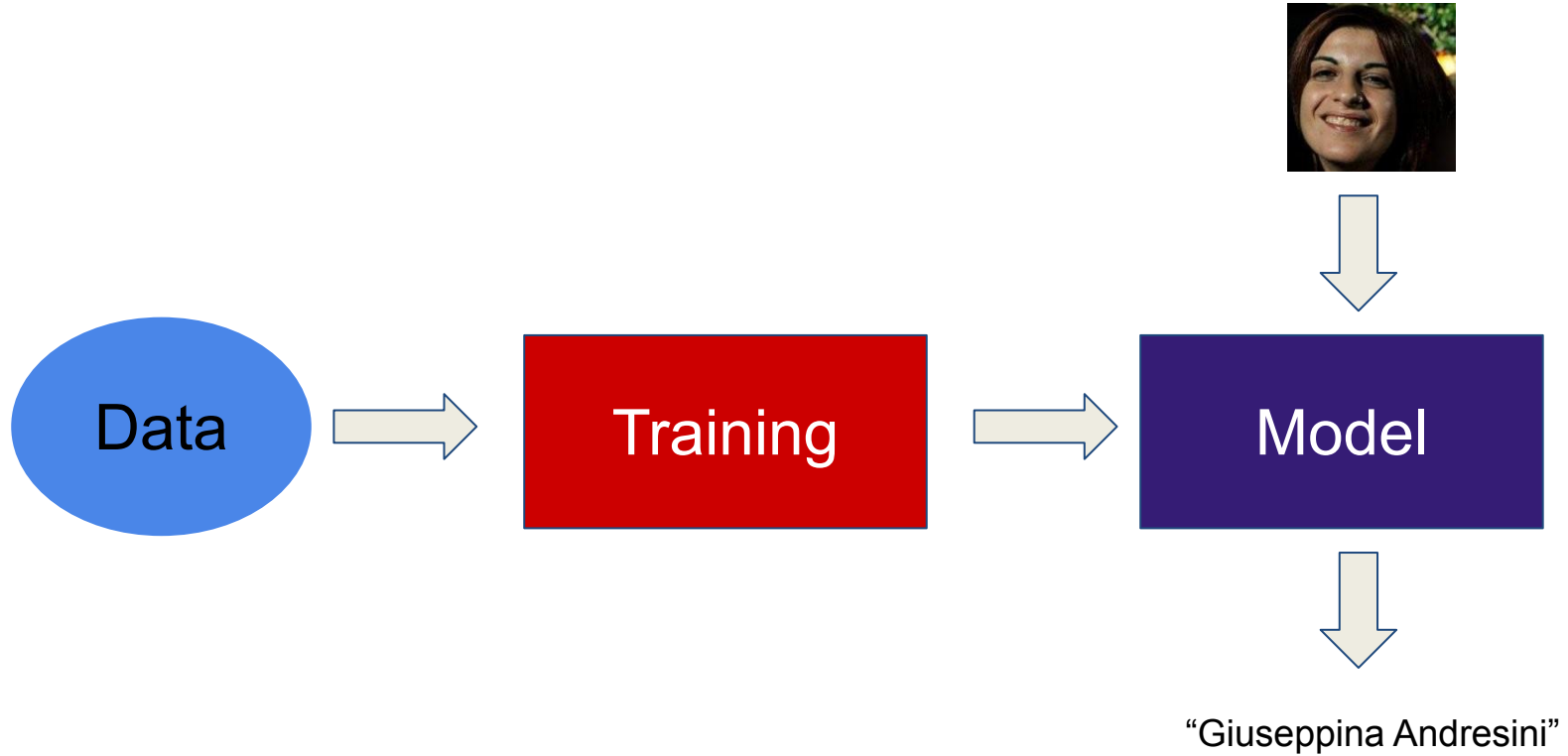


Poisoning

Strategies:

- **Data modification:**
 - Label modification of training data
 - [Biggio et al 2011 - Support Vector Machines Under Adversarial Label Noise]
- **Data injection:**
 - Augment data with adversarial examples
 - [Barreno et al 2001 - Can machine learning be secure?]
- **Logic corruption:**
 - The adversary has the ability to meddle with the learning algorithm
 - [Kloft and Laskov 2010 - Online Anomaly Detection under Adversarial Impact]

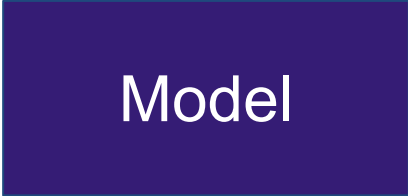
Poisoning



Poisoning: data injection

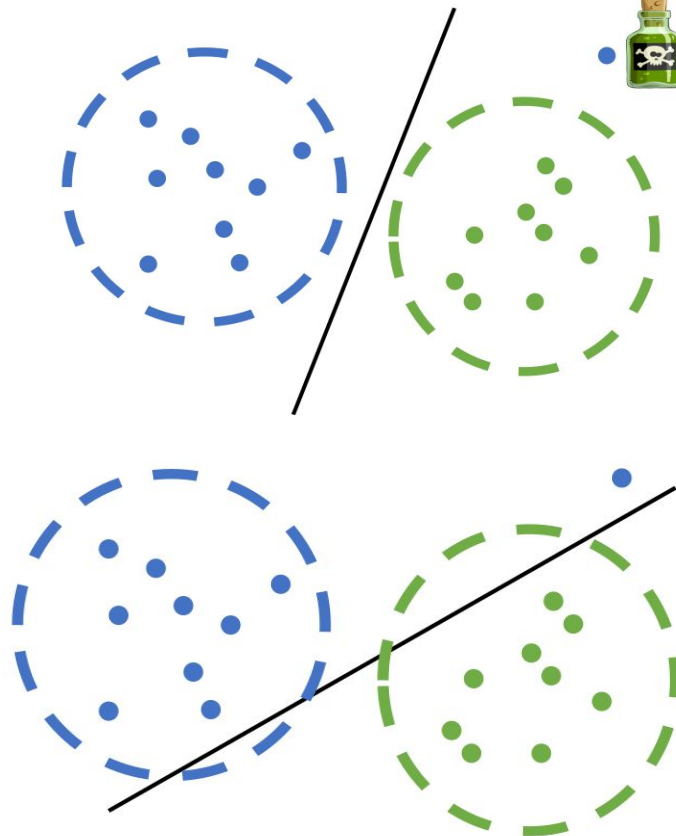
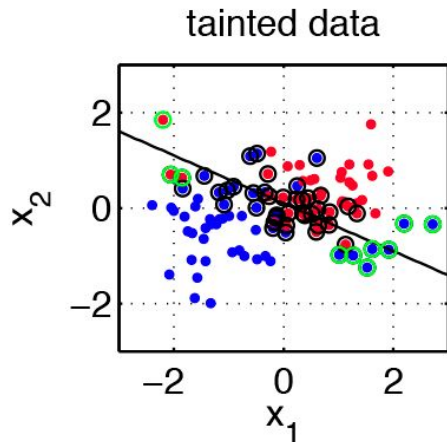
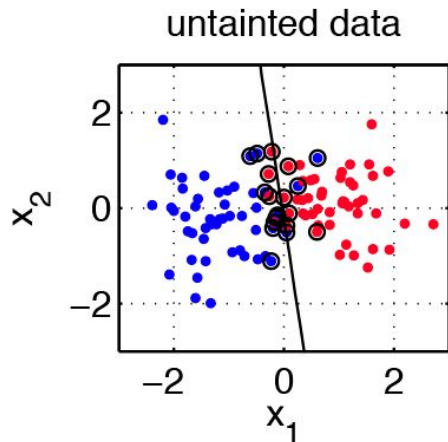


✘ 30



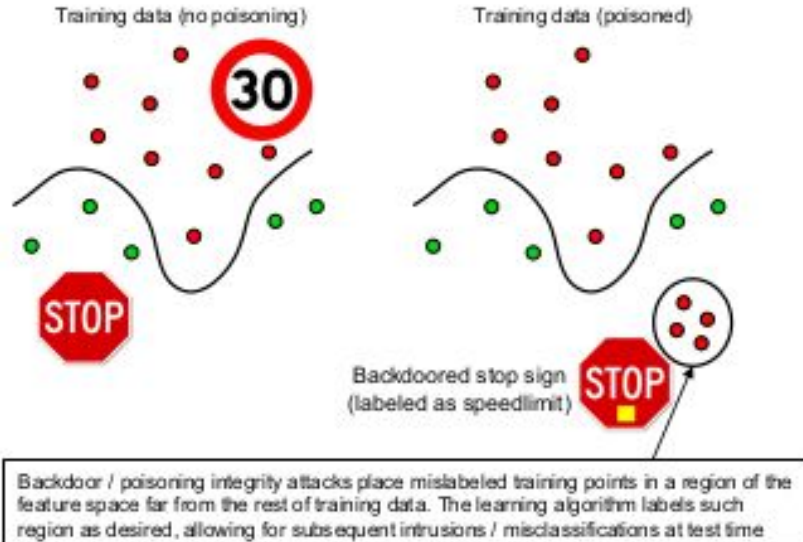
"Jennifer Lawrence"

Poisoning: data modification

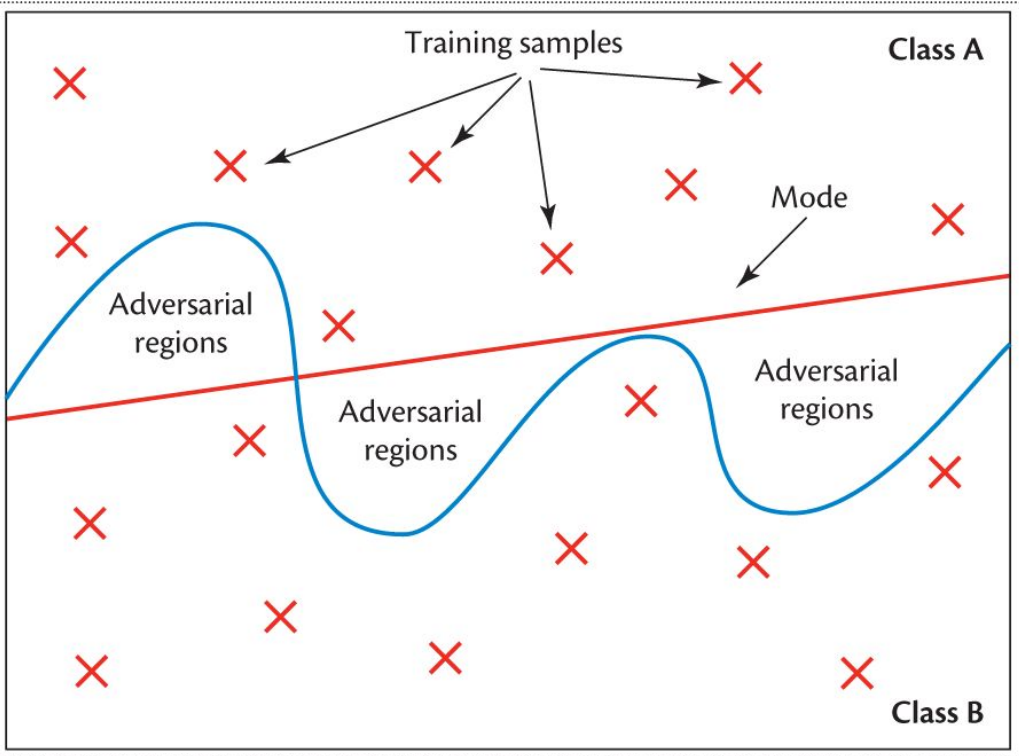
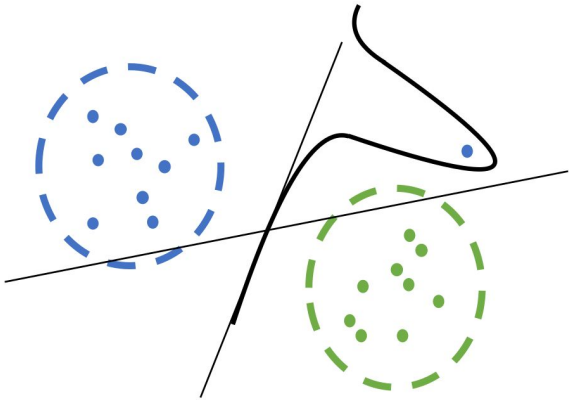


[Biggio et al (2011) -Support Vector Machines Under Adversarial Label Noise]

Poisoning: data modification

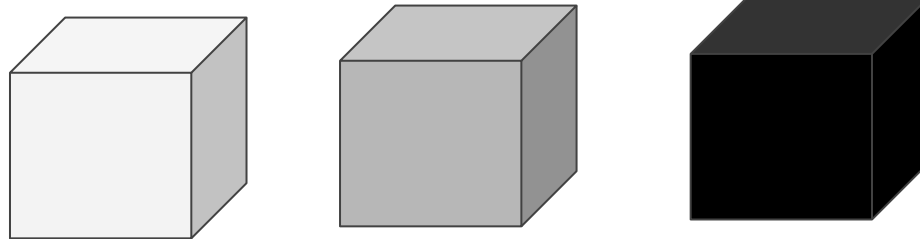


Poisoning: deep neural network



Evasion: attacker's knowledge

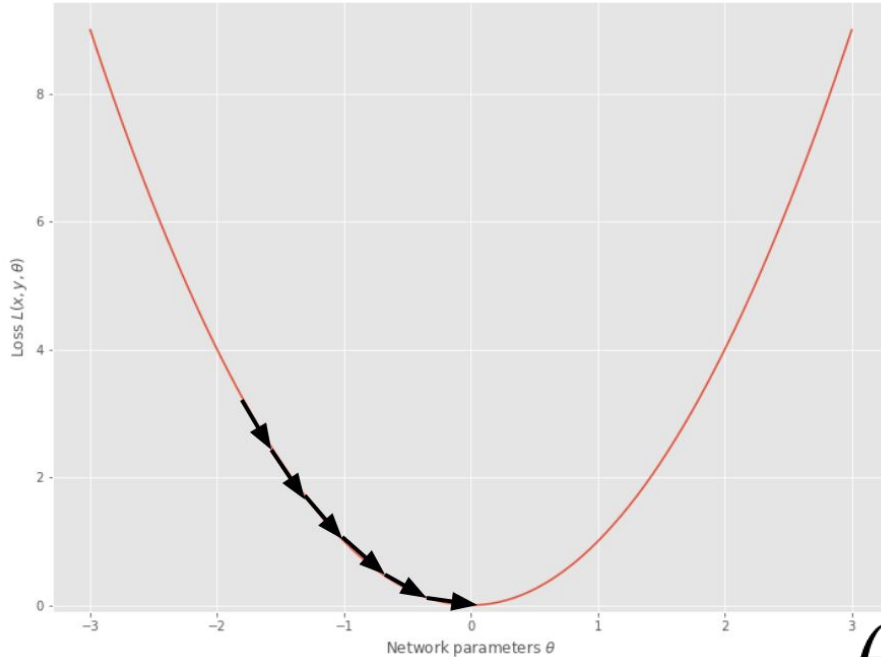
- **White box** attacks:
 - Knowledge about structure and parameters of the target model
- **Gray box** attacks:
 - Knowledge about feature representation or learning algorithm
- **Black box** attacks:
 - No knowledge about feature representation and parameters of the model



White box attacks

- Fast Gradient Sign Method (FGSM) - 2015
- Jacobian Based Saliency Map (JSMA) -2016
- DeepFool - 2016
- Carlini & Wagner - 2017
-

Training neural network



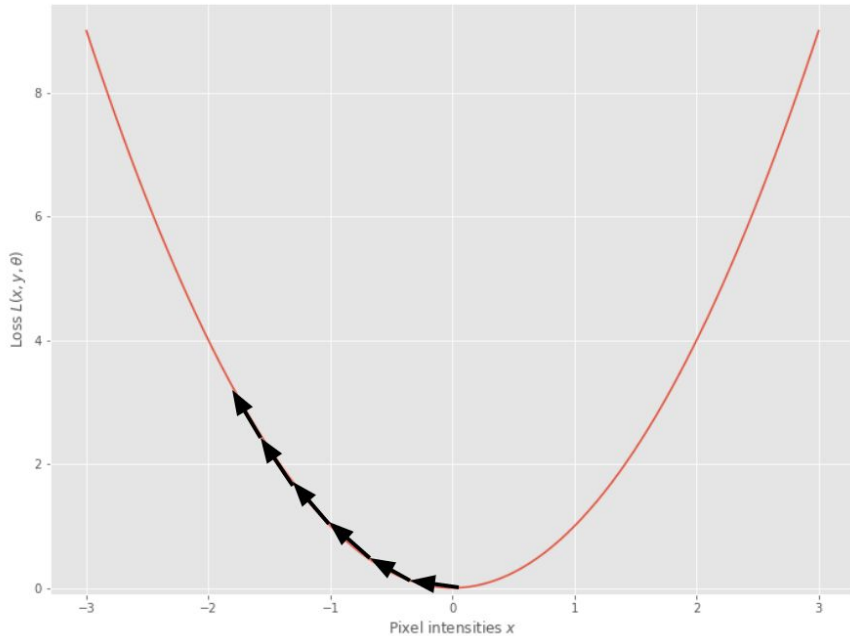
Goal: finding the **minimum** of this loss function is equivalent to finding a good set of network parameters.

$$L(\mathbf{x}, y, \theta) = (f_{\theta}(\mathbf{x}) - y)^2$$

How: differentiate this loss with respect to the parameters (theta) . Update the parameters such that the loss on that sample will decrease.

$$\theta' = \theta - \alpha \nabla_{\theta} L(\mathbf{x}, y, \theta)$$

Adversarial sample



Goal: **increase** the loss of the model on the sample x

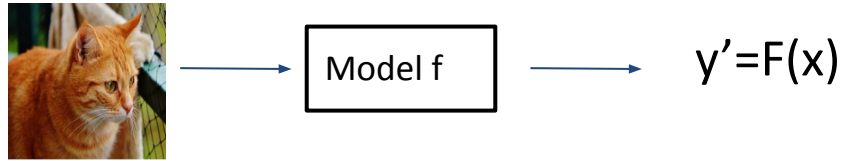
$$\mathbf{x}' = \mathbf{x} + \alpha \nabla_{\mathbf{x}} L(\mathbf{x}, y, \theta)$$

Adversarial sample

Let us consider a sample (x, y) and a model F that maps x to a predicted class label y' .

An adversarial sample x_{adv} for an original sample x is defined as:

$x_{adv} = x + \eta$ where η is a **small** perturbation added to the input such that $f(x_{adv}) \neq y$

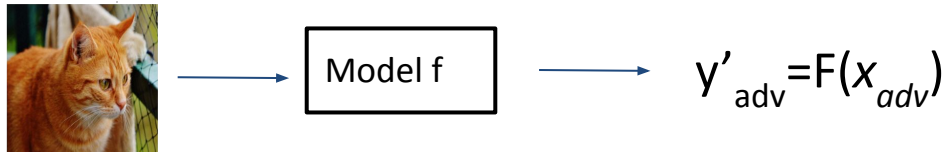


Attack goal: target or untarget

+

η

\neq



Adversarial attacks: norm

Given an original sample x an adversarial attack aims to find an adversarial example x_{adv}

$x_{adv} = x + \eta$ where η is a **small** perturbation added to the input

with $\|x - x_{adv}\| = \|\eta\| < \epsilon$ ← perturbation bound

- The constraint has the objective of disallowing perturbations which could make x unrecognisable (imperceptible perturbation)
 - without altering much the original sample, attacks can perturb a few pixels strongly (L_0), all pixels slightly (L_∞) or a mix of both (L_1 and L_2)

Adversarial attacks: norm



- ϵ (perturbation bound) is measured with a mathematical norm
 - L_0 : minimize the number of elements (pixels) modified in x_{adv} such that $x \neq x_{adv}$
 - e.g., in sticker added to stop-sign all the background is preserved only a tiny fraction of the environment is modified
 - L_1 : minimize the Manhattan distance (sum of the total perturbation values) for each pixel to create the adversarial sample.
 - $L_1 = |x^1 - x_{adv}^1| + |x^2 - x_{adv}^2| + \dots + |x^n - x_{adv}^n|$
 - Attack quite uncommon
 - L_2 : minimize Euclidean distance (MSE) for each pixel as upper bound to create adv. sample
 - $L_2 = \sqrt{(x^1 - x_{adv}^1)^2 + (x^2 - x_{adv}^2)^2 + \dots + (x^n - x_{adv}^n)^2}$
 - Commonly attack
 - L_∞ : what is the maximum value/change to *any* of the pixels in the x_{adv} image
 - $L_\infty = \max(|x^1 - x_{adv}^1| + |x^2 - x_{adv}^2| + \dots + |x^n - x_{adv}^n|)$

Evasion attacks: norm

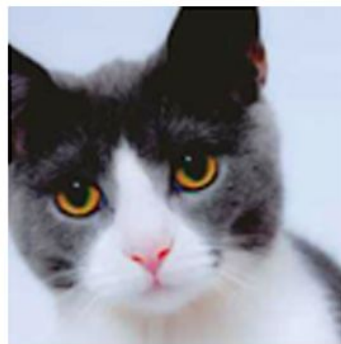


- Different meaning:
 - L_0 : a pixel that changes by 0.0001 is as influential as one that changes by 100, since the metric only account the number of pixels to change
 - L_∞ : the pixel with the maximum change is taken into account.
 - This norm is the commonly used in adversarial attacks

Adversarial attack: norm

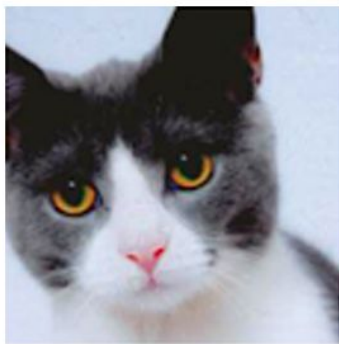
Norm: way to measuring the small change in original samples to create adversarial samples (distance from original sample)

Original



egyptian cat
(28%)

l_2 -norm=10



traffic light
(97%)

l_∞ -norm=0,05



traffic light
(96%)

l_0 -norm=5000
(sparse)



traffic light
(80%)

Adversarial attack: norm



Way to measuring the small change in original samples to create adversarial samples

	Norm bound?			
Access to gradient?	L_0	L_1	L_2	L_∞
Untarget			Carlini Wagner DeepFool	PGD FGSM BIM
Target	JSMA			

Fast Gradient Sign Method (FGSM)



“panda”

57.7% confidence

+ .007 ×



noise

=



“gibbon”

99.3% confidence

[Goodfellow et al 2015 Explaining and Harnessing Adversarial Examples]

Fast Gradient Sign Method (FGSM)

Goal: maximize **the loss** for the constructed adversarial sample

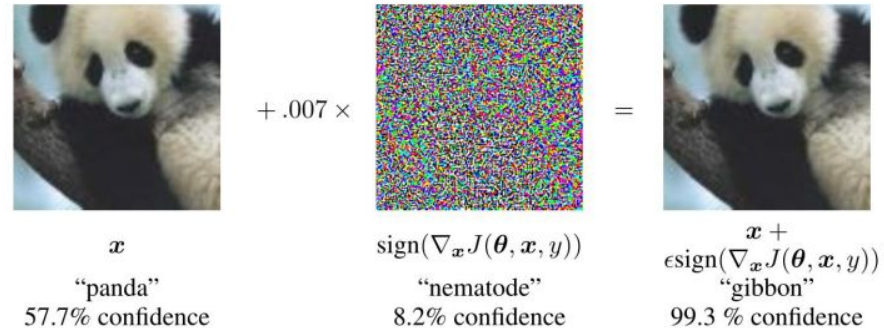
$x_{adv} = x + \eta$ with respect to the true label y

$$\max L(F(x+\eta), y)$$

instead of optimizing the parameters to decrease loss (with images constant) we optimize the image pixels to increase loss (constant parameters)

with $\|\eta\|_{\infty} < \epsilon$

- Single step: perturbation is computed once by following the gradient of the function F
- Constraint optimization: the attack applies a bound equal to L_{∞}



Fast Gradient Sign Method (FGSM)

$$\text{Adversarial example: } X_{\text{adv}} = X + \eta$$

$$\eta = \epsilon \text{ sign}(\nabla_x L(x, y, W))$$

Perturbation

Gradient loss function with respect to the image x for the true label y

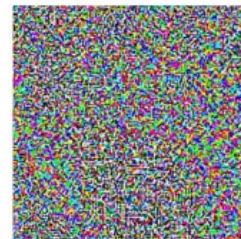
*small scalar
(e.g., 0.1, 0.007)*

distort pixels in the opposite direction of the loss with respect to the target class y



x
"panda"
57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$
"nematode"
8.2% confidence

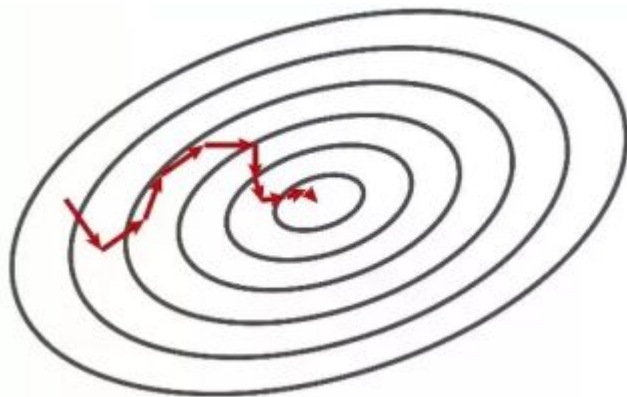
=



$x + \epsilon \text{ sign}(\nabla_x J(\theta, x, y))$
"gibbon"
99.3 % confidence

Iterative multistep algorithms

- FGSM is a **single-step gradient update** : limits the power of attack
- More powerful: iterative object optimization running steepest descent for multiple iterations
- Two extension of FGSM:
 - **BIM**
 - **PGD**



Basic Iterative Method (BIM)

- BIM is an extension of FGSM
 - Add noise in multiple iterations with a step size α
 - After each iteration the result is clipped to ensure that the perturbation is within ϵ -neighbourhood of original sample x (maximum perturbation for each pixel)

$$x_{\text{adv}}^0 = x$$

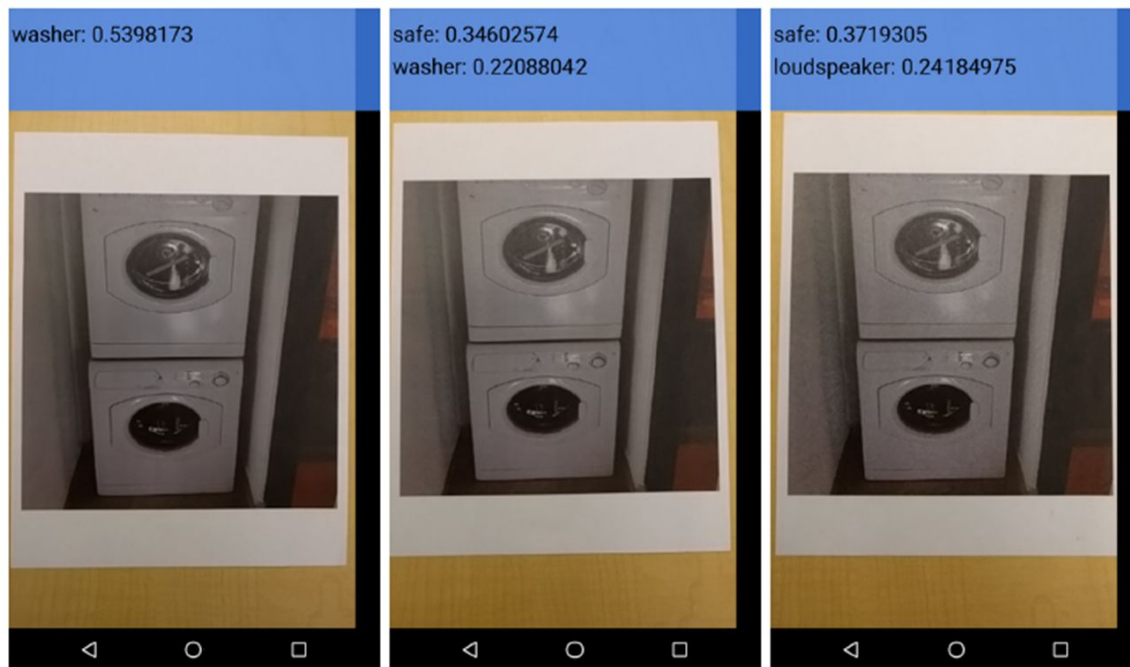
$$x_{\text{adv}}^t = \text{Clip}_{x,\epsilon} \{ x_{\text{adv}}^{t-1} + \alpha \cdot \text{sign}(\nabla_x L(f(x_{\text{adv}}^{t-1}), y)) \}$$

a = step size of iteration

$$\text{Clip}_{x,\epsilon} \{ x_{\text{adv}} \} = \min \{ 255, x + \epsilon, \max \{ 0, x - \epsilon, x_{\text{adv}} \} \}$$

t = number of iterations

Basic Iterative Method (BIM)



(b) Clean image

(c) Adv. image, $\epsilon = 4$

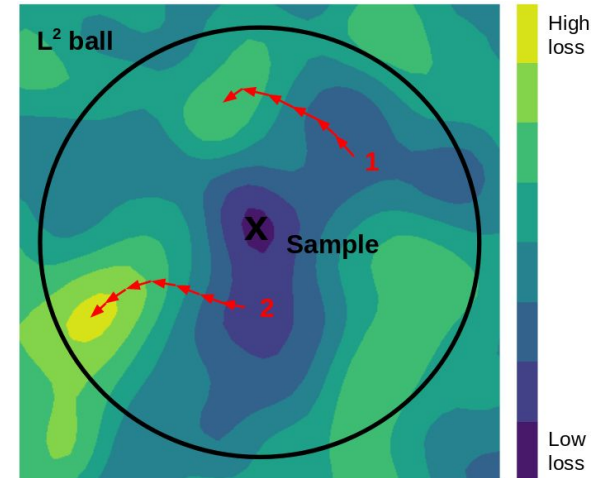
(d) Adv. image, $\epsilon = 8$

Project Gradient Descent (PGD)

PGD is an extension of BIM , where after each step of perturbation, the adversarial example is projected back onto the ϵ -ball of x (decided by norm) using a projection function Π

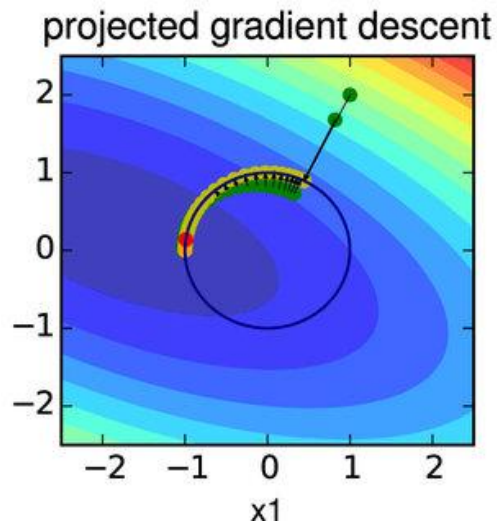
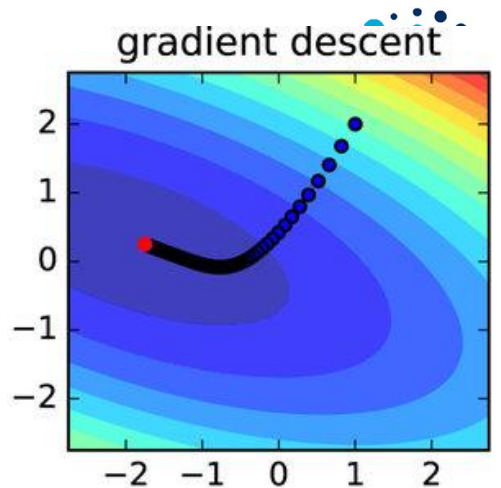
$$x_{\text{adv}}^t = \Pi_{\epsilon} (x^{t-1} + \alpha \cdot \text{sign}(\nabla_x L(f(x^{t-1}), y)))$$

α =gradient
step size



Project Gradient Descent (PGD)

- **Gradient Descent:** is a standard way to solve unconstrained optimization problem
 - $\min_{x \in \mathbb{R}} f(x) \rightarrow$ any x in \mathbb{R} can be a solution
- **Project Gradient Descent:** is a standard way to solve constrained optimization problem.
 - $\min_{x \in Q} f(x) \rightarrow$ not any x in \mathbb{R} can be a solution but inside the set Q
 - given a point x , PGD try to find a point in Q which is closest to x



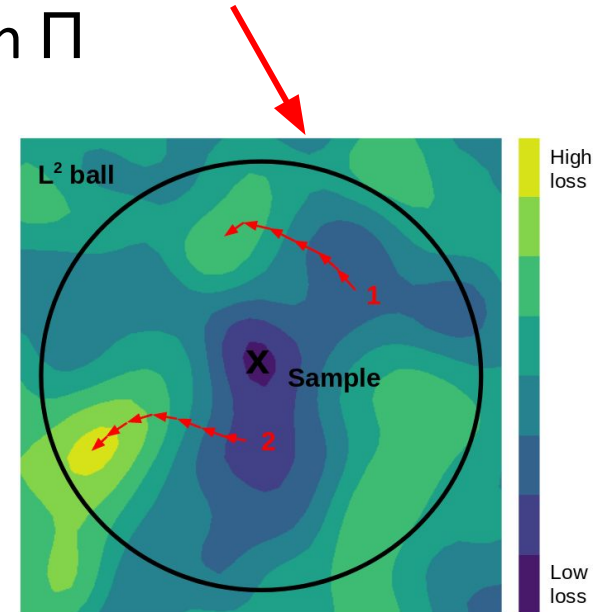
Project Gradient Descent (PGD)

PGD is an extension of BIM , where after each step of perturbation, the adversarial example is **projected back onto the ϵ -ball of x** (decided by L norm) using a projection function Π

$$x_{adv}^t = \Pi_{\epsilon} (x^{t-1} + \alpha \cdot \text{sign}(\nabla_x L(h(x^{t-1}), y)))$$

Different from BIM, PGD uses random initialization for x for each iteration by adding random noise from a uniform distribution with values in the range $(-\epsilon, \epsilon)$

[Madry et al 2019 Towards Deep Learning Models Resistant to Adversarial Attacks]

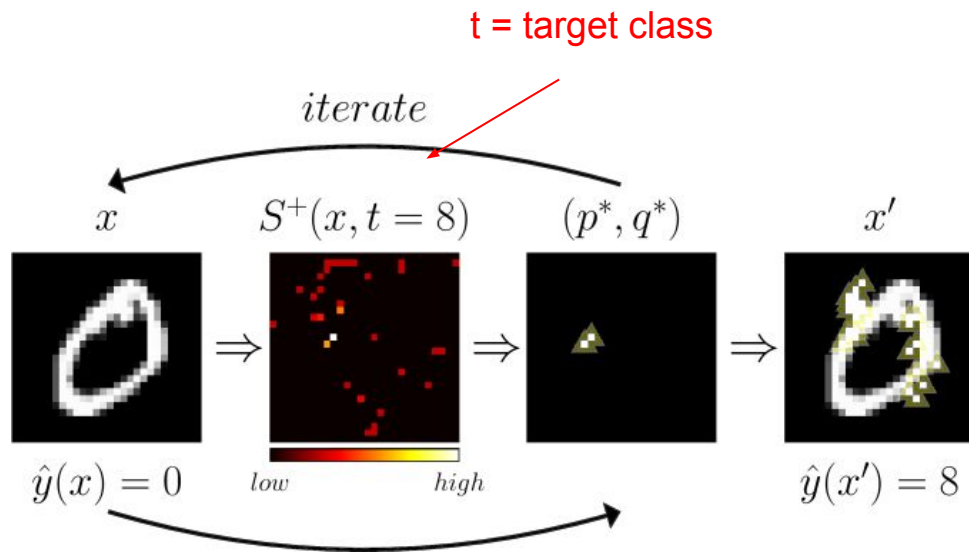


Project Gradient Descent (PGD)



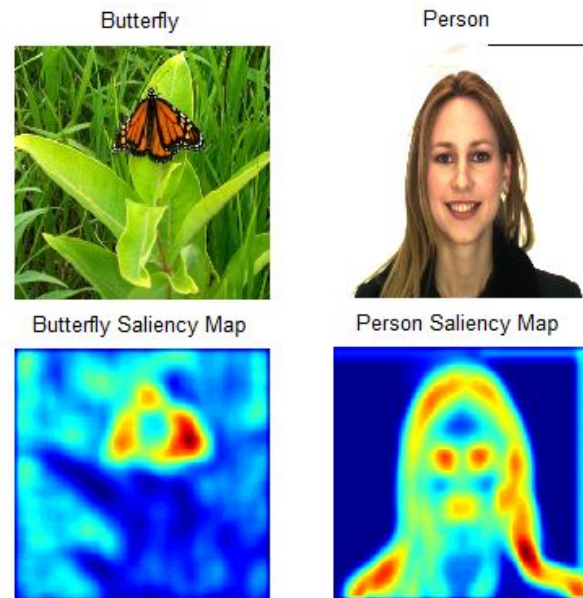
Jacobian-based Saliency Map Attack (JSMA)

- JSMA is a target attack -> goal to create an adversarial sample that is misclassified in a specified class
- JSMA different from FGSM that alter all pixel uses L_0 norm that greedily modifies pairs of pixels at a time.
 - using a saliency map, which shows an impact each pixel has on the classification result.
 - A large value means, that changing this pixel will have a significant impact on the outcome of the classification.



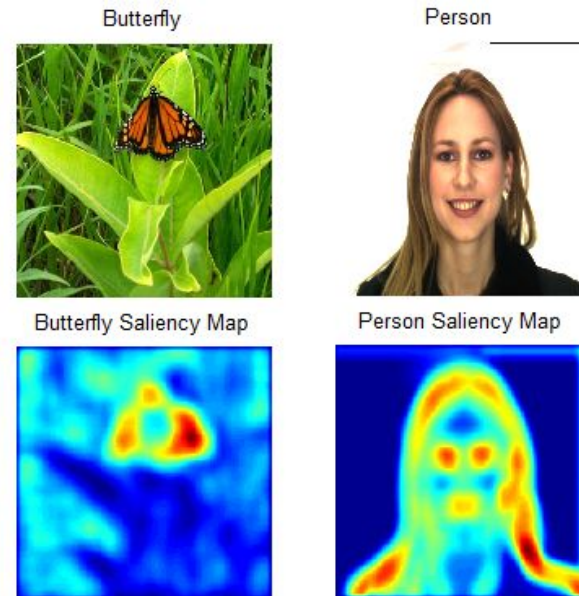
JSMA: saliency map

- Saliency Map: visual explanation of the predictions of a classifier
 - gradient of each input feature x^i (e.g., each pixel) to the class score (i.e., how influential to predict a particular class c)
 - $y'(x) = \arg \max_c f(x)(c)$ where $f(x)$ is the softmax probabilities vector predicted by the model.



JSMA: adversarial saliency map

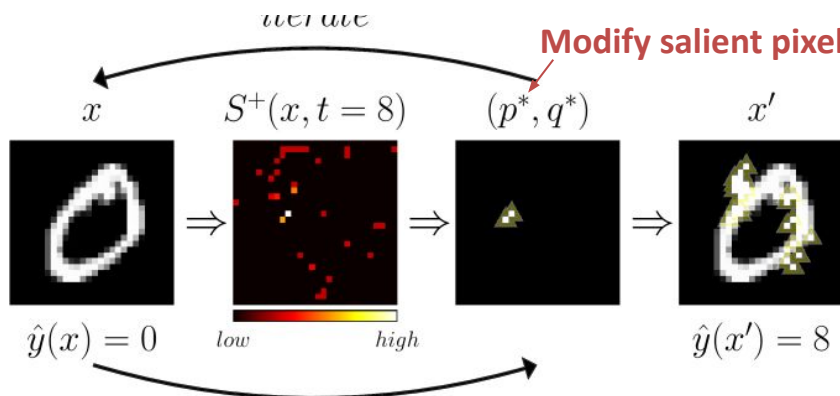
- Adversarial saliency map: visual explanation of which input features should perturb in order to predict the adversarial sample as a target class.
- Adversary want to misclassify a sample x that is assigned to a target t
 - $f_t(x)$ must be increased while the confidence for $f_j(x)$ for all other classes $j \neq t$ decrease until $t = \arg \max f(x)$



$S(x, t)$ = measure how a feature (x^i) is positively correlates with t , while also negatively correlates with all other classes $j \neq t$

rejects input with a negative derivative or an overall positive derivative on other classes.

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



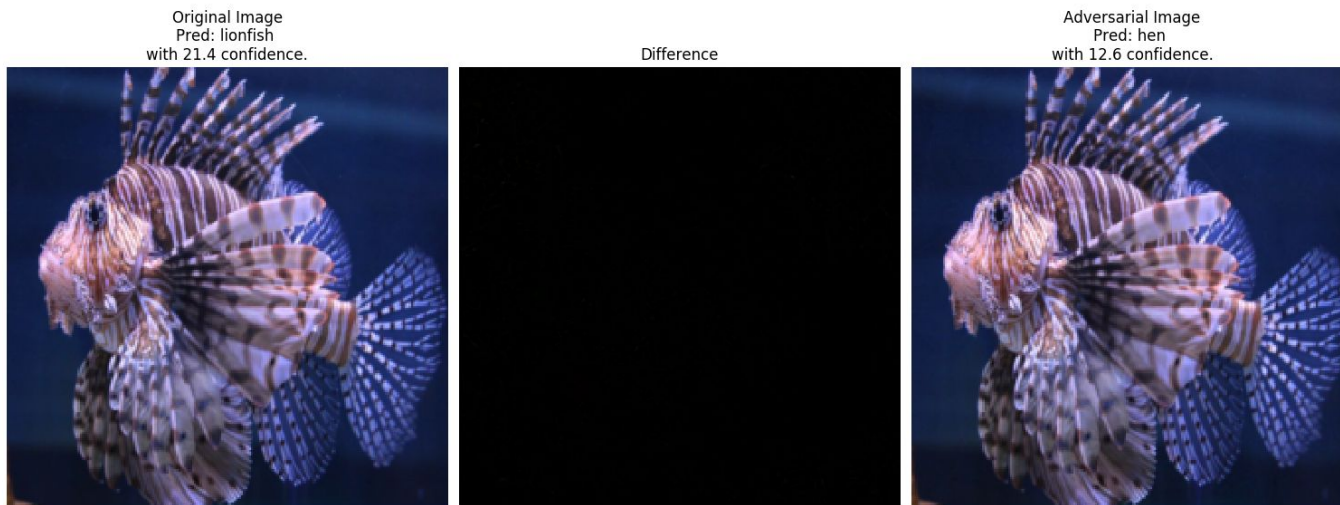
How much $F_t(x)$ will increase and $F_j(x)$ will decrease given a modification of the input feature i

DeepFool

Compute a minimal norm adversarial perturbation for a given image in iterative manner,
to find the decision boundary closest to the normal sample

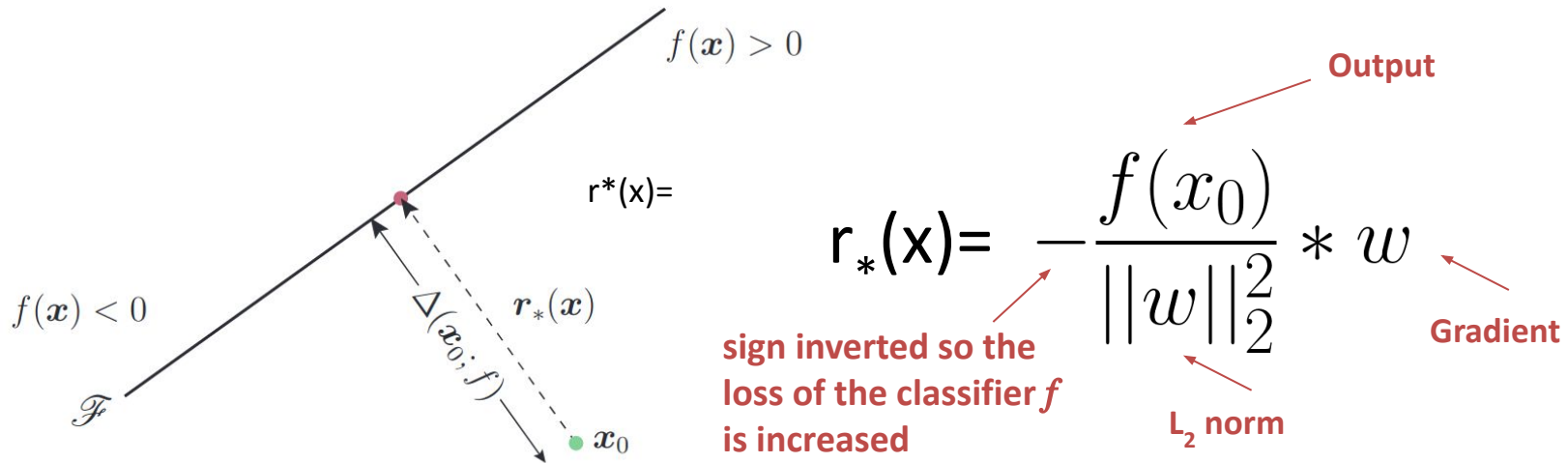
Perform steps by linearly approximate the decision function, according to the l_2 norm

Perturbation is small since gradient is orthogonal to the boundary



DeepFool for binary classifier

- Using a linear binary classifier, that the **robustness** of a model (f) for an input x_0 is equal to the distance of x_0 to the hyperparameter plane (which separates the 2 classes)
- The minimal perturbation $r^*(x)$ to change the classifier's decision must project the input image of x_0 orthogonal to the hyperplane of classification



DeepFool



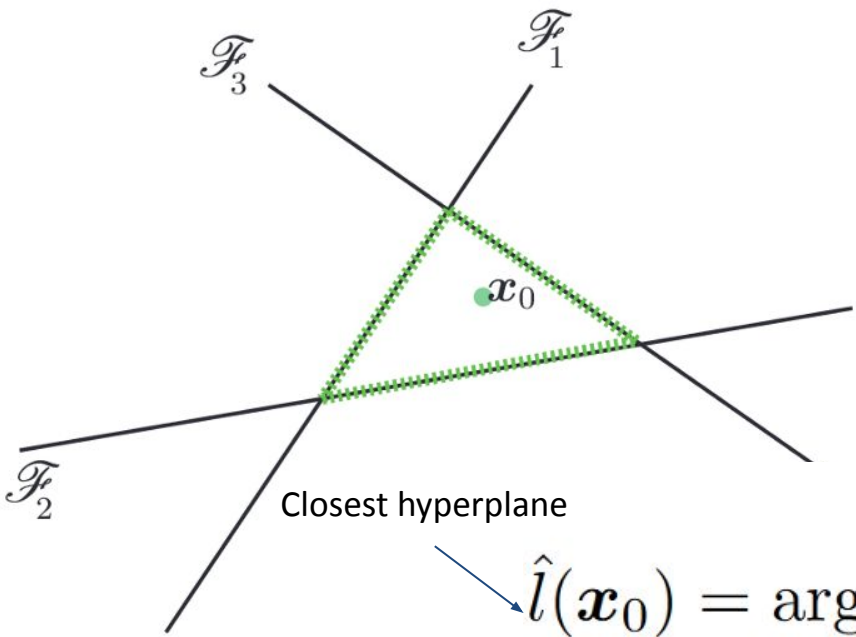
- **!!!Misclassification is not guarantee since models are not linear in nature**
 - To alleviate this issue, the algorithm works iteratively and adds the previous perturbation to the next perturbation, which is performed until the label changes or max iterations are reached

Algorithm 1 DeepFool for binary classifiers

- 1: **input:** Image \mathbf{x} , classifier f .
 - 2: **output:** Perturbation $\hat{\mathbf{r}}$.
 - 3: Initialize $\mathbf{x}_0 \leftarrow \mathbf{x}, i \leftarrow 0$.
 - 4: **while** $\text{sign}(f(\mathbf{x}_i)) = \text{sign}(f(\mathbf{x}_0))$ **do**
 1. Start and continue loop while the true label and the label of the adversarially perturbed image is the same.
 - 5: $\mathbf{r}_i \leftarrow -\frac{f(\mathbf{x}_i)}{\|\nabla f(\mathbf{x}_i)\|_2} \nabla f(\mathbf{x}_i)$,
 2. Calculate the projection of the input onto the closest hyperplane (minimal perturbation)
 - 6: $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \mathbf{r}_i$,
 3. Add that perturbation to the image and test
 - 7: $i \leftarrow i + 1$.
 - 8: **end while**
 - 9: **return** $\hat{\mathbf{r}} = \sum_i \mathbf{r}_i$.
-

DeepFool for multiclass

Multiple binary classifiers



- The minimum perturbation needed would be from the closest hyperplane to x_0 .
- Given there are multiple classes, the loss and backpropagation would need to be computed **for each class label** allowed by the function.
- When finding the minimum perturbation, the difference must be taken between the computed values for each label and the computed values for the label of the original prediction.

$$\hat{l}(\mathbf{x}_0) = \arg \min_{k \neq \hat{k}(\mathbf{x}_0)} \frac{|f_k(\mathbf{x}_0) - f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_0)|}{\|\mathbf{w}_k - \mathbf{w}_{\hat{k}(\mathbf{x}_0)}\|_2}$$

most probability after the true class

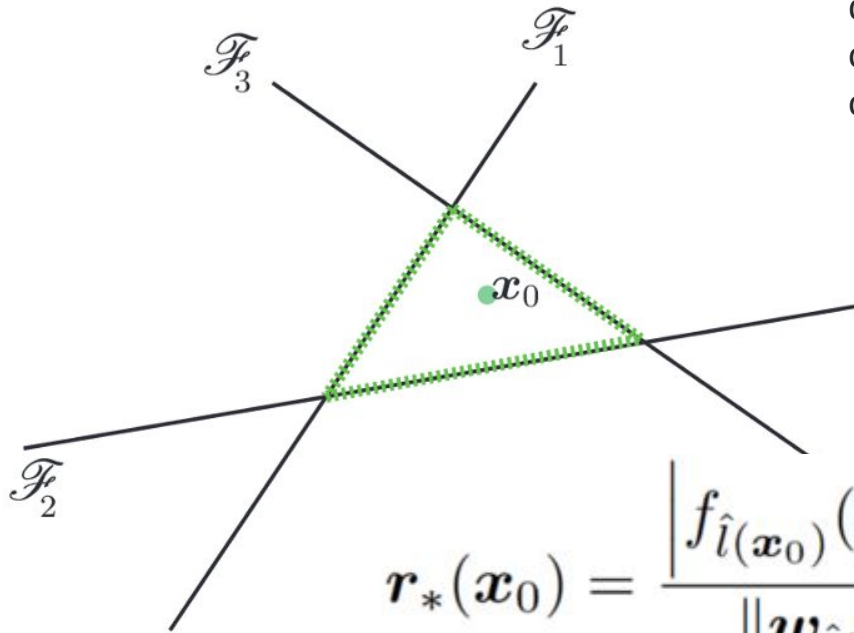
true class

DeepFool for multiclass



Multiple binary classifiers

- The minimum perturbation $r_*(x_n)$ is the vector that project differences of the classifier output and the gradients of the outputs, as well as taking the absolute value of the model output.



$$r_*(x_0) = \frac{|f_{\hat{l}(x_0)}(x_0) - f_{\hat{k}(x_0)}(x_0)|}{\|w_{\hat{l}(x_0)} - w_{\hat{k}(x_0)}\|_2} (w_{\hat{l}(x_0)} - w_{\hat{k}(x_0)})$$

DeepFool for multiclass



Algorithm 2 DeepFool: multi-class case

- 1: **input:** Image \mathbf{x} , classifier f .
 - 2: **output:** Perturbation $\hat{\mathbf{r}}$.
 - 3:
 - 4: Initialize $\mathbf{x}_0 \leftarrow \mathbf{x}$, $i \leftarrow 0$.
 - 5: **while** $\hat{k}(\mathbf{x}_i) = \hat{k}(\mathbf{x}_0)$ **do**
 - 6: **for** $k \neq \hat{k}(\mathbf{x}_0)$ **do**
 - 7: $\mathbf{w}'_k \leftarrow \nabla f_k(\mathbf{x}_i) - \nabla f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_i)$ ← 1. Store the difference between the original gradients and the gradients of each of classes (w_k) and the difference between labels (f_k)
 - 8: $f'_k \leftarrow f_k(\mathbf{x}_i) - f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_i)$
 - 9: **end for**
 - 10: $\hat{l} \leftarrow \arg \min_{k \neq \hat{k}(\mathbf{x}_0)} \frac{|f'_k|}{\|\mathbf{w}'_k\|_2}$ ← 2. Calculate the minimal vector that projects c on the closes hyperplane
 - 11: $\mathbf{r}_i \leftarrow \frac{|f'_i|}{\|\mathbf{w}'_i\|_2} \mathbf{w}'_i$
 - 12: $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \mathbf{r}_i$
 - 13: $i \leftarrow i + 1$
 - 14: **end while**
 - 15: **return** $\hat{\mathbf{r}} = \sum_i \mathbf{r}_i$ ← 3. Total perturbation is the sum over all calculated perturbation
-

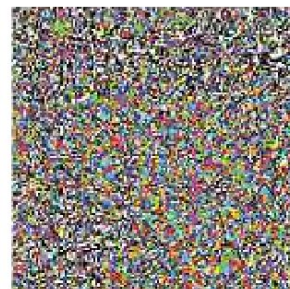
DeepFool vs FGSM



True label: whale



Predict: turtle
DeepFool



Predict: turtle
FGSM

Carlini Wagner (CW) attack

- Goal: find a small change η that make an image x misclassified but so that the result is still a valid image/example
 - All the other attacks were using constraints on perturbation (ϵ), here the score is used as a penalizer, modulated by C .
- How close we are getting to being classified as t

$$f = [1 - C(x + \eta)_t]$$

$C(x + \eta)_t = t$ is satisfied if $f(x + \eta) \leq 0$ is satisfied

Probability of $x + \eta$ to be classified as t .

If the probability is low value of f is closer to 1 value of f is closer to 1 whereas when it is classified as t , f is equal to 0.

Carlini Wagner (CW) attack

- How close we are getting to being classified as t

$$f(x') = \max (\max\{Z(x')_i\} - Z(x')_t, -k)$$

with $i \neq t$

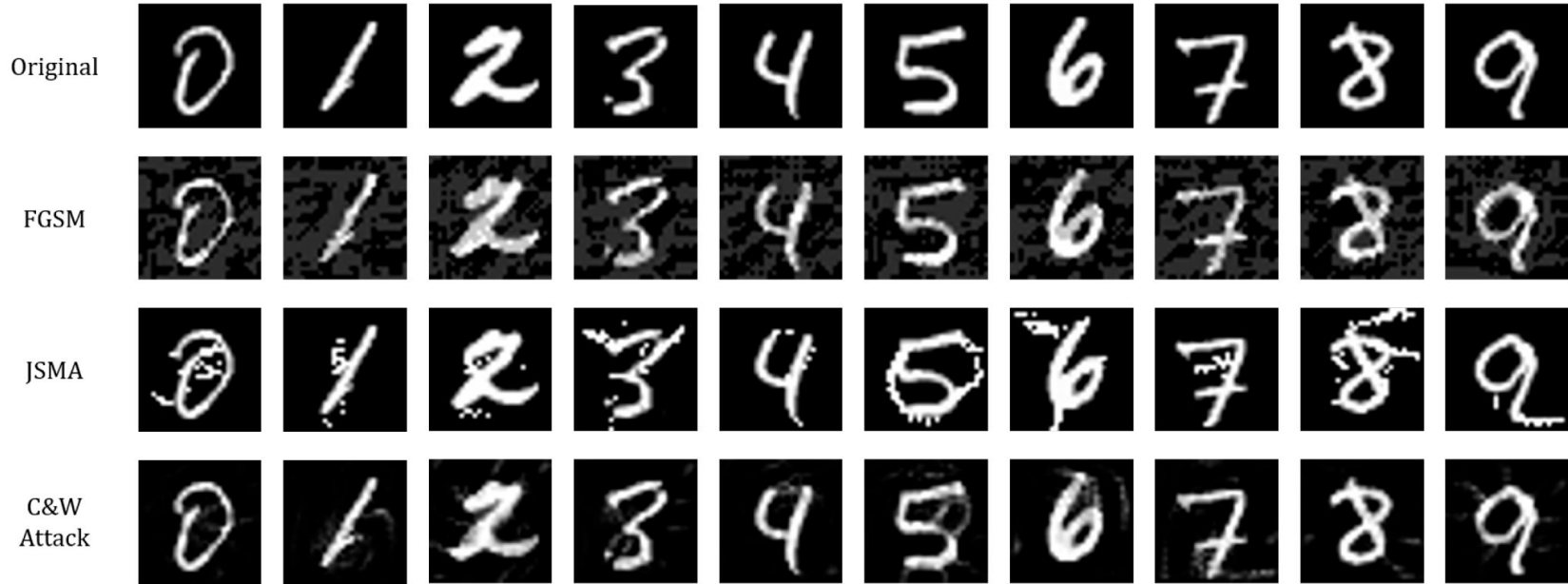
highest probabilities for non target classes

vector of probabilities for target class

lower limit of loss, at last the value in $-k$ will always hold

difference between “what the model thinks the current image most probably is” and “what we want it to think/ misclassified target”. So when the model thinks that this image is what we want it to think, this value is negative (the probability of target class is higher than any of the non target classes)

Comparison



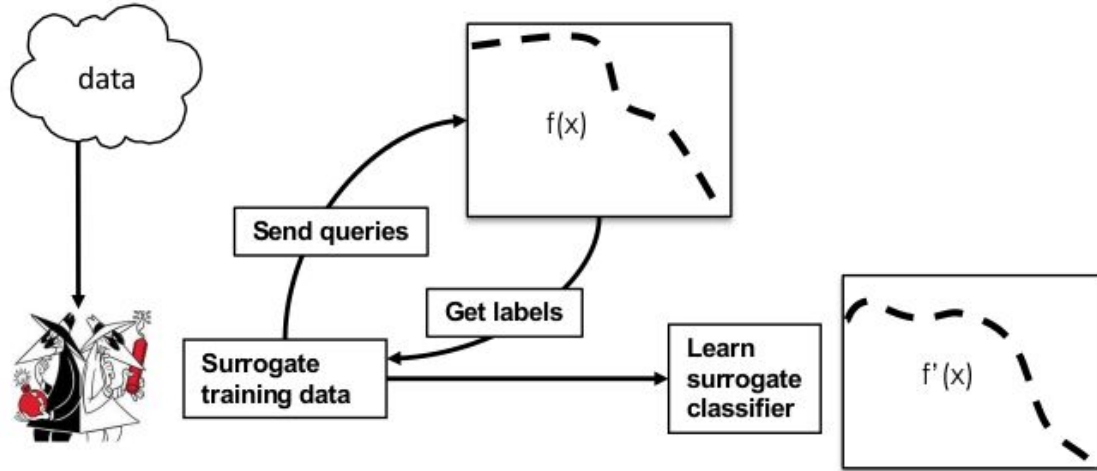
Evasion: black box



- **Transferability** based black-box attack :
 - train a substitute model, and craft adversarial examples against the substitute, and transfer them to a victim model
- It is very likely that an adversarial example of one network can fool another network
- Transferability depends on the type of attack
 - e.g. examples built with FGSM are highly transferable

Evasion: Black box

Training a local model to substitute for the target DNN, using inputs synthetically generated by an adversary and labeled by the target DNN.

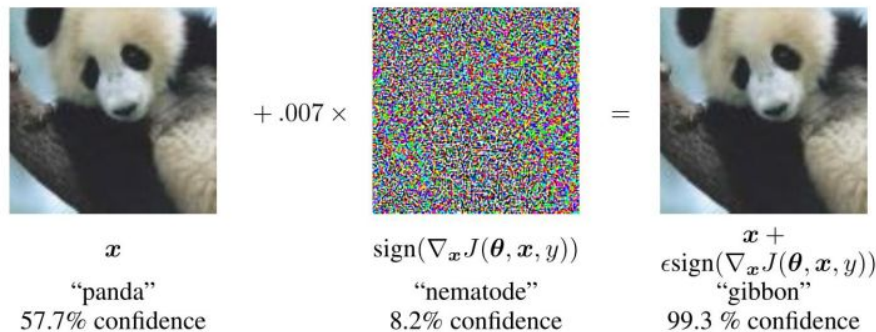


Evasion: black box

- train a substitute network based on the input/output pairs of the target network
- build adversarial examples for the substitute network
- attack the target network with the examples built for the substitute network
- due to transferability the attack is very likely to succeed

Adversarial transferability

White box scenario:
access to the model



Transferability captures the ability of an attack against a machine-learning model to be effective against a different, potentially unknown, model that was learned for the same problem (this was observed in the context of deep learning, as well as for other learning paradigms)

Adversarial sample



Target model

gibbon

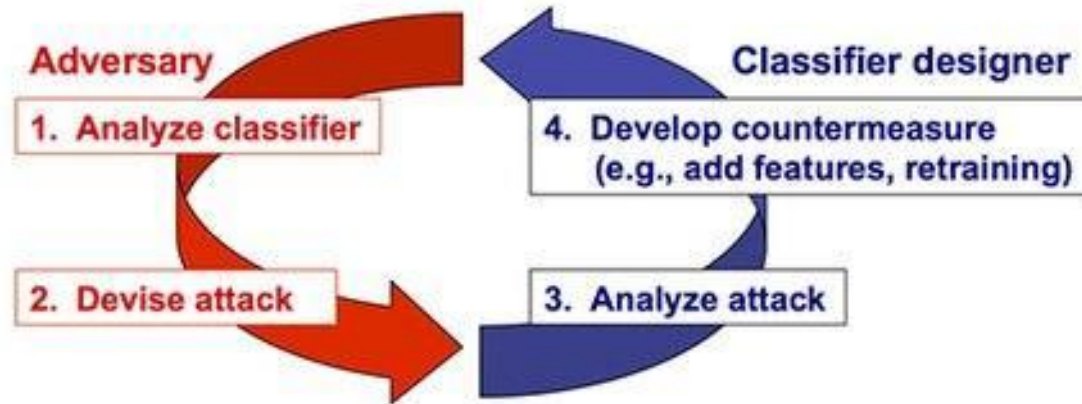
Defense Strategy

- Reactive defense
- Proactive defense



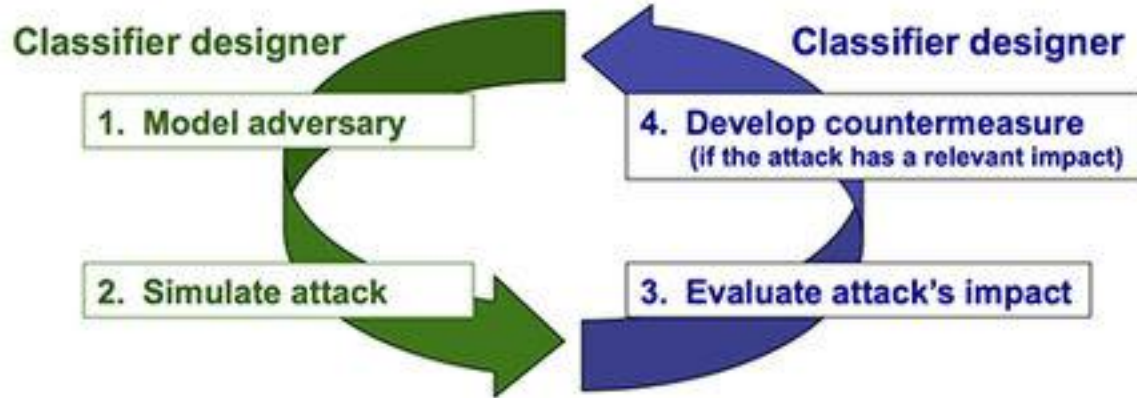
Reactive defense

- Timely detection of novel attack
- Frequent classifier retraining
- Verification of consistency of classifier against training data and ground-truth label



Proactive defense: simulating attack

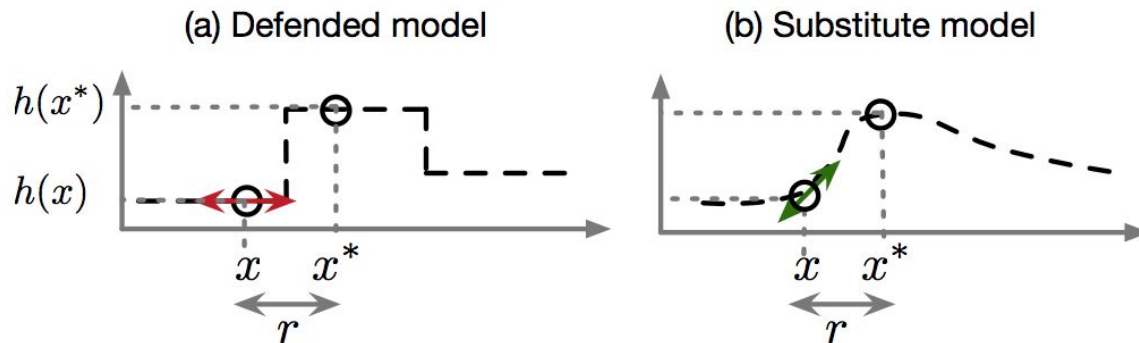
- Gradient hiding
- Adversarial training



Proactive defense: gradient hiding

- **Gradient masking:**

- reduce the sensitivity of models to small changes by finding adversarial direction using a substitute
- smooths the model's decision surface in adversarial directions
- defensive distillation

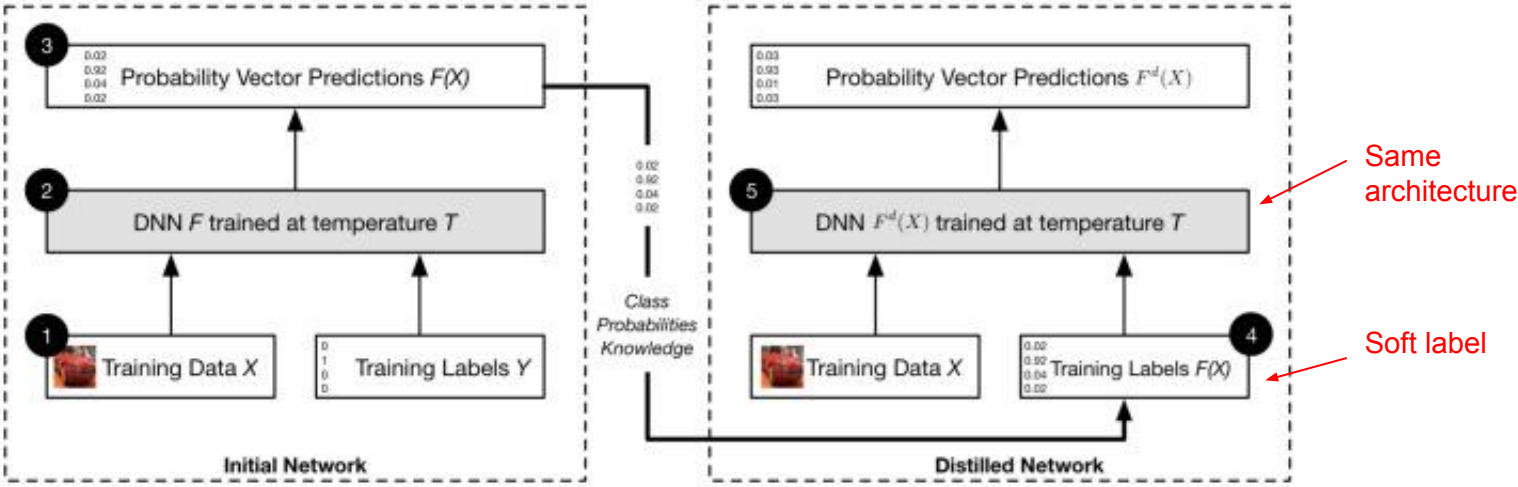


Distillation vs Defensive distillation

- Distillation was first introduced by Hinton et al. in, where the goal was for a small model to mimic a large, computationally expensive model.
 - Teacher and student
- **Defensive distillation** has a different goal:
 - two model that are the same but one model is trained to predict the probabilities output by another model that was trained earlier
 - smooths the model's decision surface in adversarial directions exploited by the adversary.
 - The second *distilled* model is more robust to attacks such as the fast gradient sign method or the Jacobian-based saliency map approach

Defensive Distillation

Training using the probability distribution as the target, not just the argmax class label

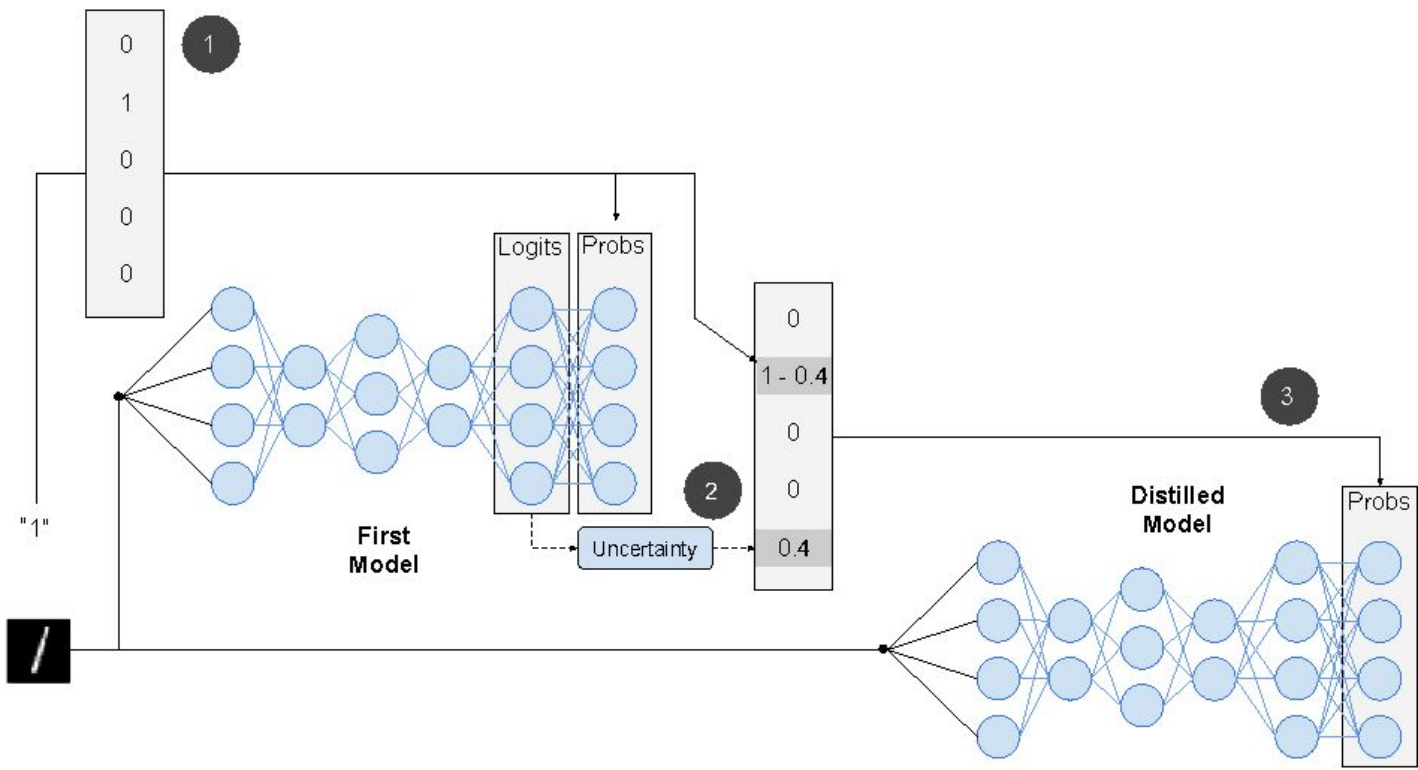


[Papernot et al 2016- Distillation as a Defense to Adversarial Perturbations against Deep Neural Networks]

Defensive distillation is not robust

- Carlini et al. (2016) prove that with a slight modification to the standard adversarial attacks, distilled networks can be attacked
- Defensive distillation work on attacks that use gradient to approximate each pixel's importance to modify to create adversarial sample (e.g., JSMA)
- Carlini et al. propose attacks using the logits to create adversarial sample (not softmax)
 - instead of taking the gradient of the inputs to the softmax, they take the gradient of the actual output of the network
 - For $T = 100$, attacks had 96.4% success rate while only changing on average 36.4 pixels; works on all T from 1 to 100

Extending defensive distillation



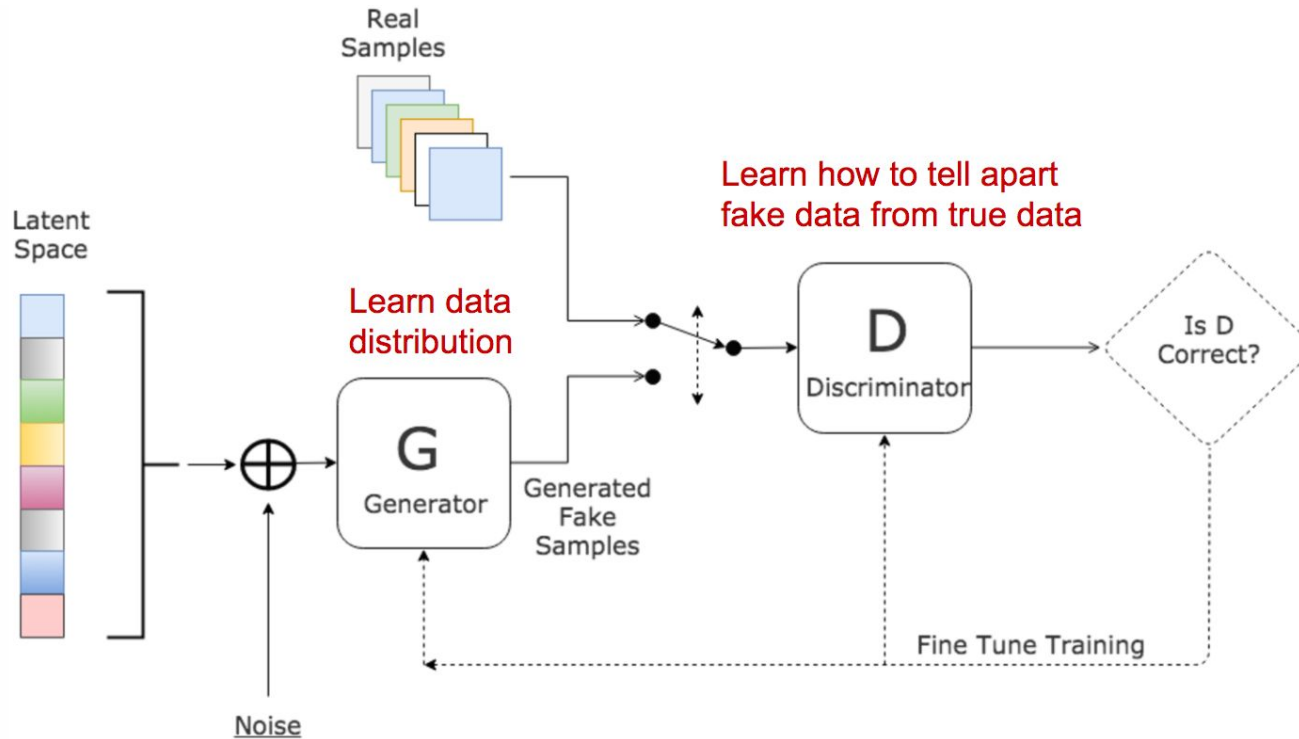
Proactive defense: adversarial training

- Introduce adversarial samples in training set to improve the robustness of the target model
 - [Szegedy et al 2013 - Intriguing properties of neural networks: inject adversarial samples and modified its labels]
 - Malik
- Punish misclassified adversarial image
 - [Huang et al 2016 - Learning with a Strong Adversary]

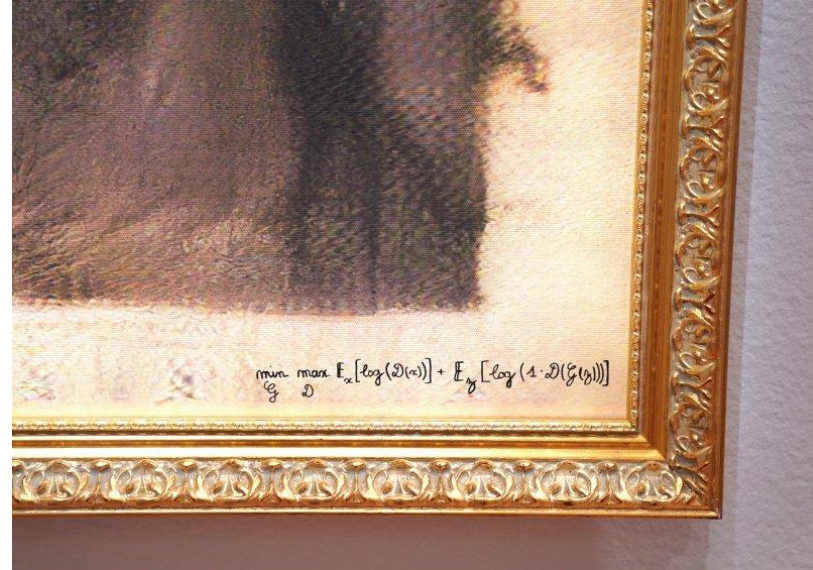
Adversarial training: data augmentation

- GAN (generative adversarial network)
- Adversarial autoencoder

GAN



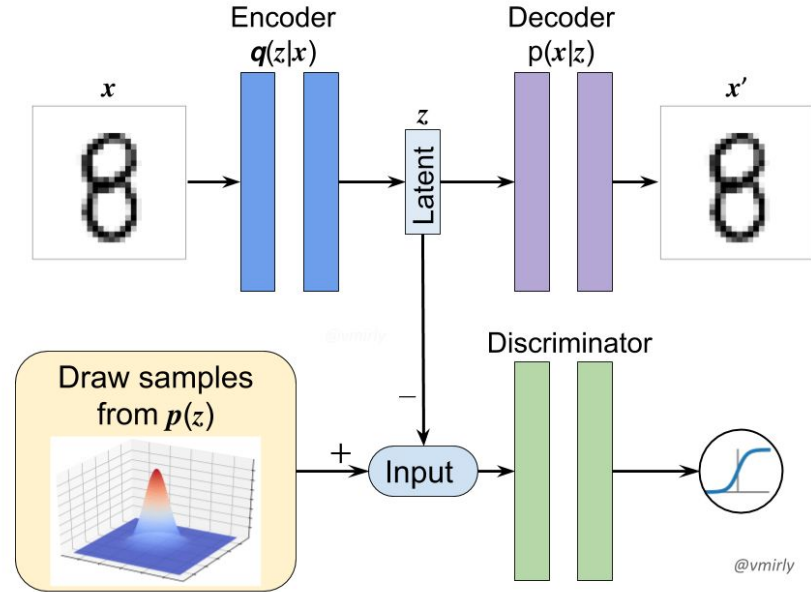
GAN



Painting above was made by a GAN and sold for **\$432 thousand**

<https://time.com/5435683/artificial-intelligence-painting-christies/>

Adversarial autoencoder



Adversarial training: feature robustness

Training set



frog

Restrict to features
of robust model



New training set



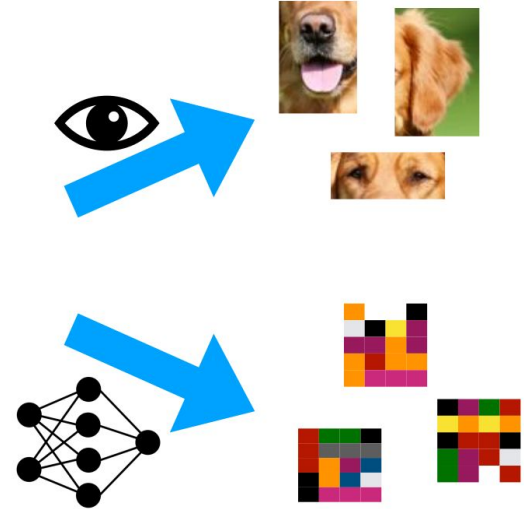
"Robustified" frog

Adversarial training: feature robustness

- Goal: robust ML as humans classification
- Classifiers tend to use any available signal to do so, even those that look incomprehensible to humans to generalize
- Adversarial examples can be attributed to the presence of non-robust features:
 - highly predictive
 - incomprehensible to humans.



dog



Adversarial training: feature robustness

Standard training: use all of features, maximize accuracy

Adversarial training: use only single robust features (at the expense of accuracy)

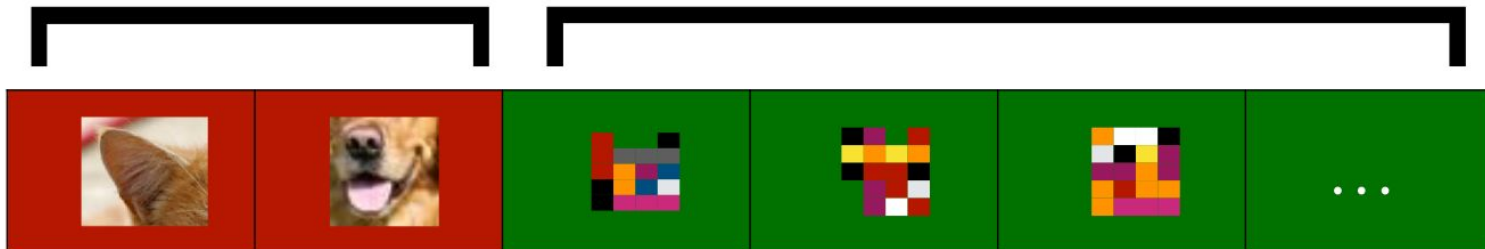
Under adversarial perturbation

Robust features

Correlated with label
even with adversary

Non-robust features

Correlated with label on average,
but can be flipped within, e.g., ℓ_2 ball



XAI for feature robustness

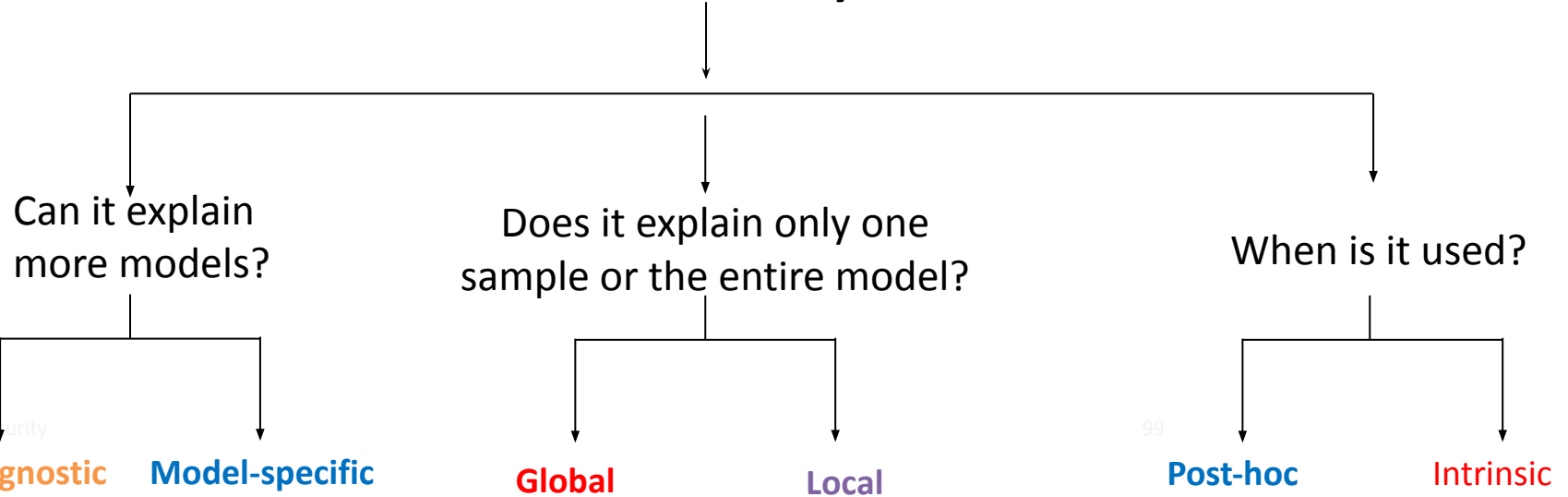


- AI-based systems (e.g., Deep Learning-based systems) are accurate, but they usually learn **black-box models**
- Unjustifiable and misleading predictions may make systems **vulnerable to attacks**
 - by leading to unsecured critical systems
- Explainable and interpretative results make AI solutions more **robust and trustworthy**
 - **Transparency** of model decisions is mandatory to:
 - provide justifiable decision making
 - produce accurate explanations of decision models' behavior
 - help to design effective countermeasures

XAI taxonomy



XAI taxonomy



XAI for cybersecurity

99

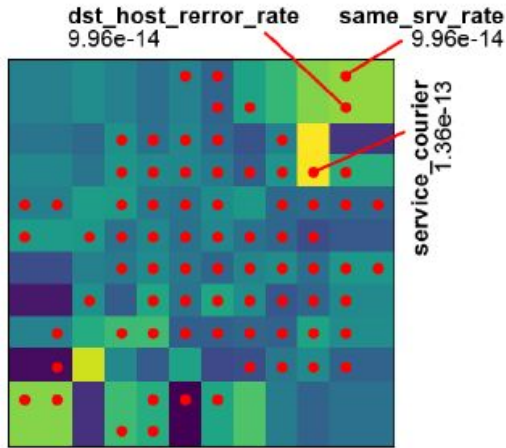
Explanation



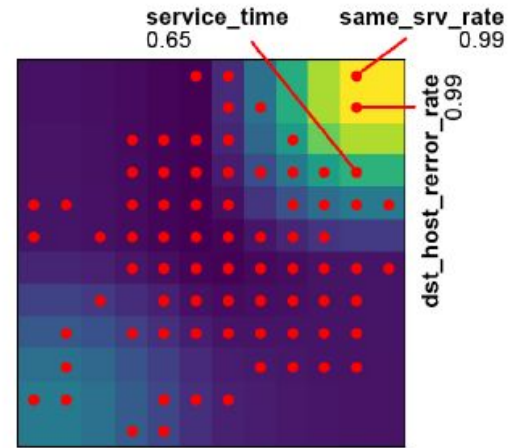
XAI for cybersecurity

Explanation of network flow attacks

NSL-KDD dataset



Top-1 Grad-CAM
(attack)



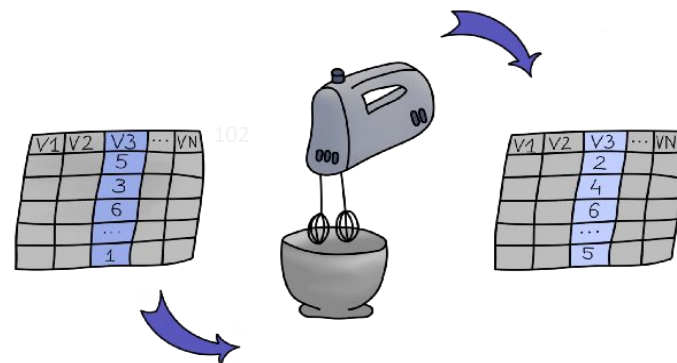
Top-2 Grad-CAM
(attack)

Dalex

(moDEL Agnostic Language for Exploration and eXplanation)

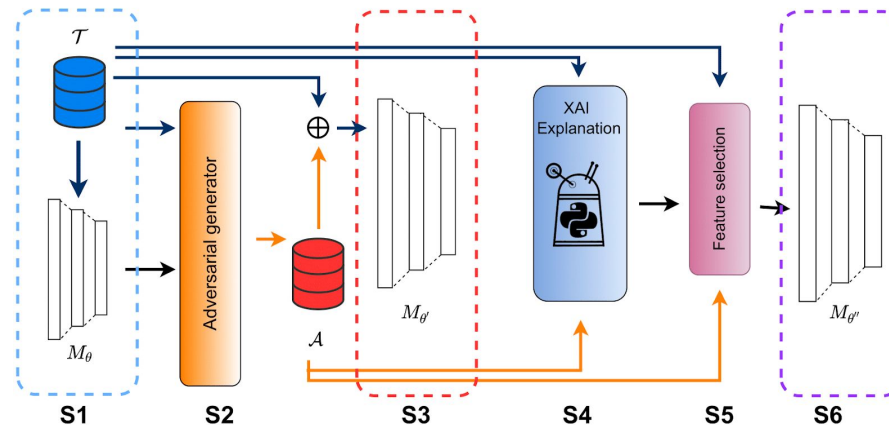


- **Model agnostic** and **post-hoc** technique
- To understand both the **global** and **local** structure of black box models
- To explain the behavior of models by measuring the global relevance of features on decisions
 - It uses permutation-based feature-importance
 - It permutes the value of each feature
 - It computes a loss function before and after the permutation
 - If the loss increases then the feature is important



Robust feature selection with XAI

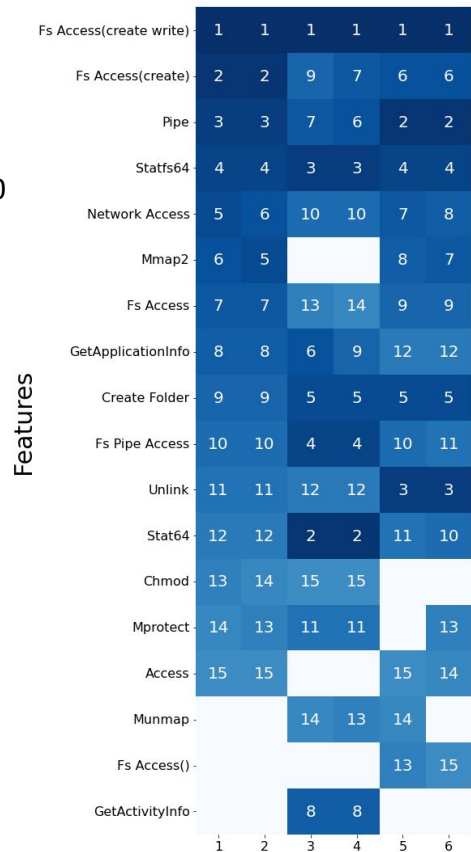
- Adversarial training to learn robust deep neural models
- Post-hoc global explanations with **DALEX**
 - To perform feature selection by extracting the top-k features ranked by DALEX on the training set



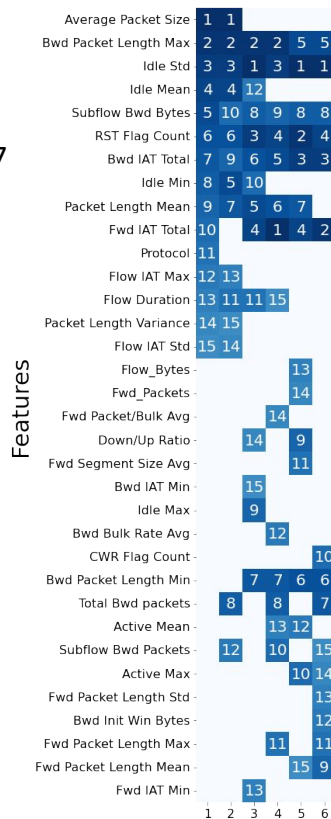
Robust feature selection with DALEX



MalDroid20



CICIDS17



Legend

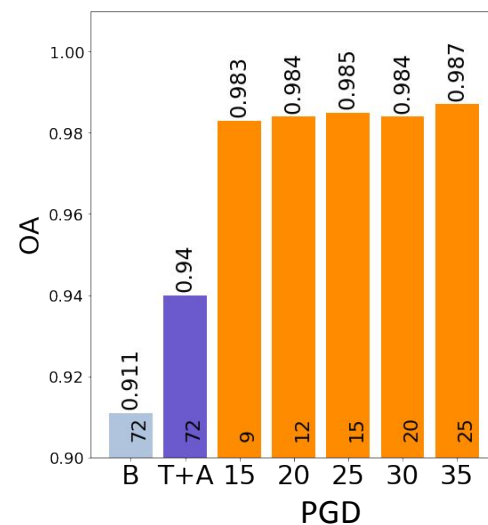
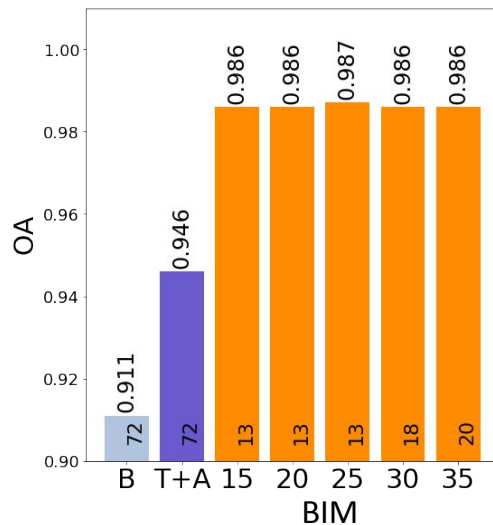
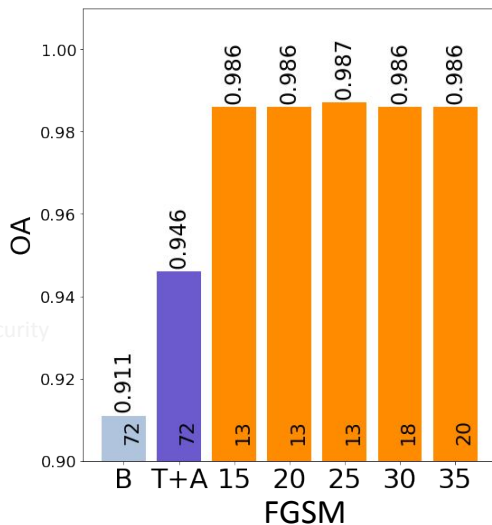
- (1) T+A (FGSM) explains T
- (2) T+A (FGSM) explains A
- (3) T+A (BIM) explains T
- (4) T+A (BIM) explains A
- (5) T+A (PGD) explains T
- (6) T+A (PGD) explains A

104

Robust feature selection with DALEX

- **B**: which discards both adversarial training and XAI-based feature selection
- **T+A**: which discards XAI-based feature selection
- **T+A+XAIFS (n)**: which implements both adversarial training and XAI-based feature selection (with n number of feature selected)

CICIDS2017

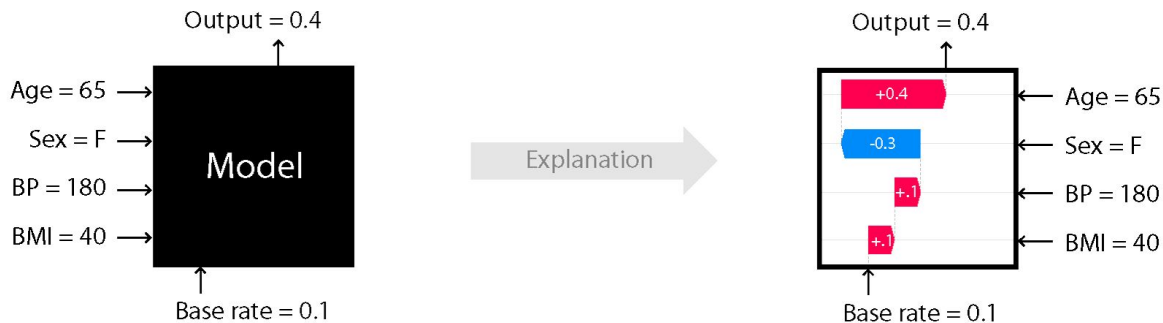


SHAP

(SHapley Additive exPlanation)

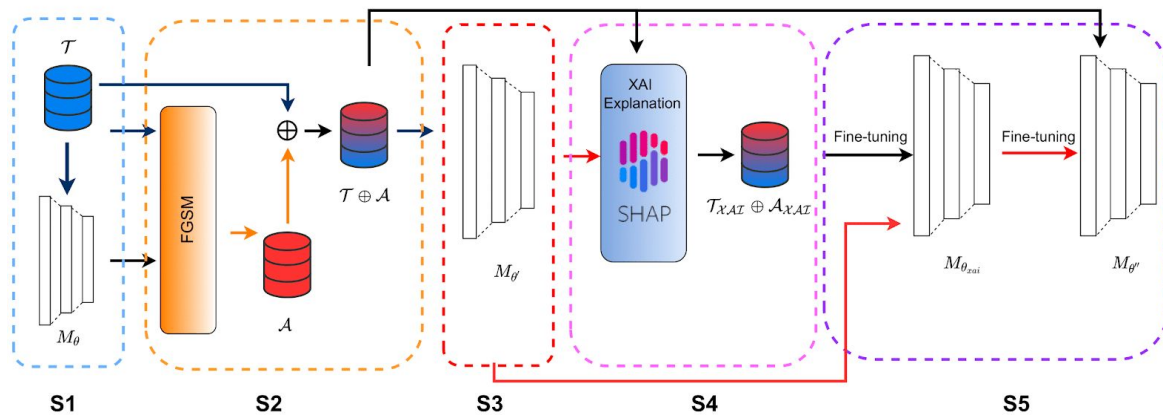


- **Model agnostic** and **post-hoc** technique
- It produces **local** explanations of decisions for single samples
- Based on the game theory, Shapley Values
 - SHAP values measure the average contribution of each feature on a decision and show whether a feature has a positive or negative effect on a decision



Improving accuracy with SHAP

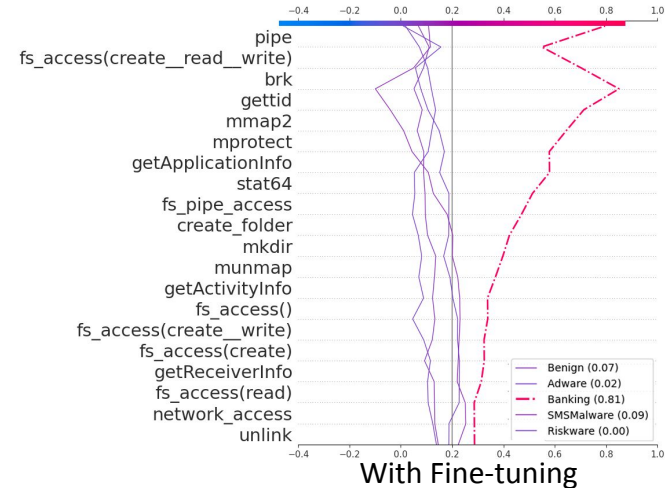
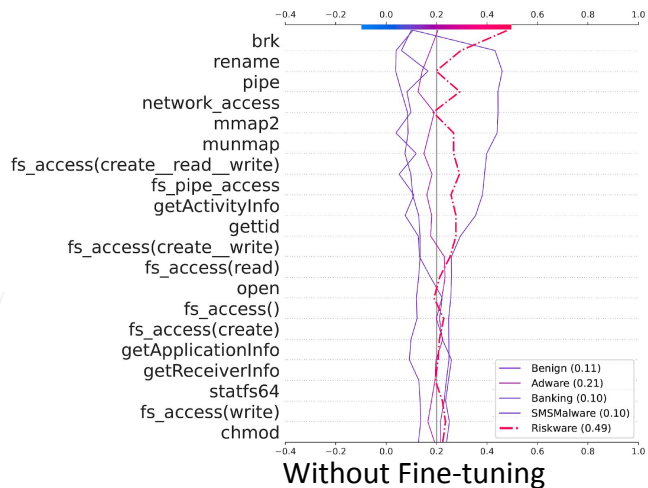
- Combination of Adversarial Training and XAI
 - Adversarial training with FGSM to create a new deep neural model robust to adversarial samples
 - SHAP to extract the local feature importance for each training sample
 - Fine-tuning of the the deep neural model driven by XAI values



XAI for cybersecurity

Improving accuracy with SHAP

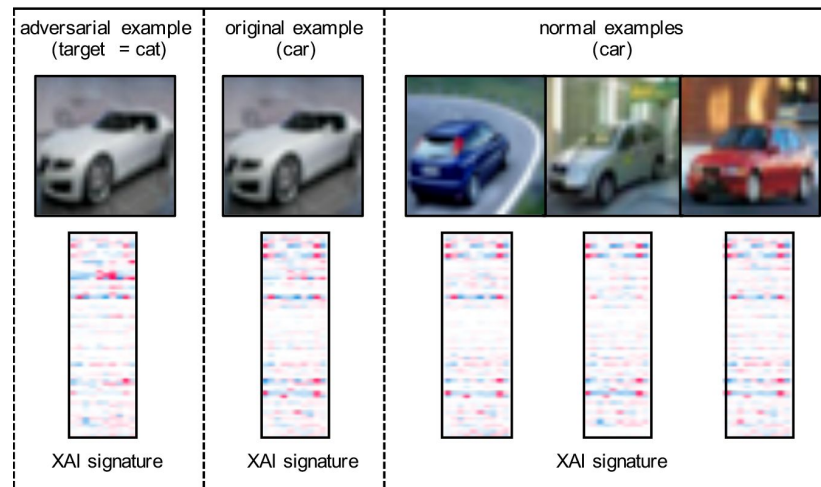
- Analysis of SHAP explanations for the class Banking (testing samples) on Maldroid20 dataset
 - A decision yielded with a model trained using Adversarial training and XAI-based fine-tuning is better separated from a decision yielded with a model trained without fine tuning



XAI for cy

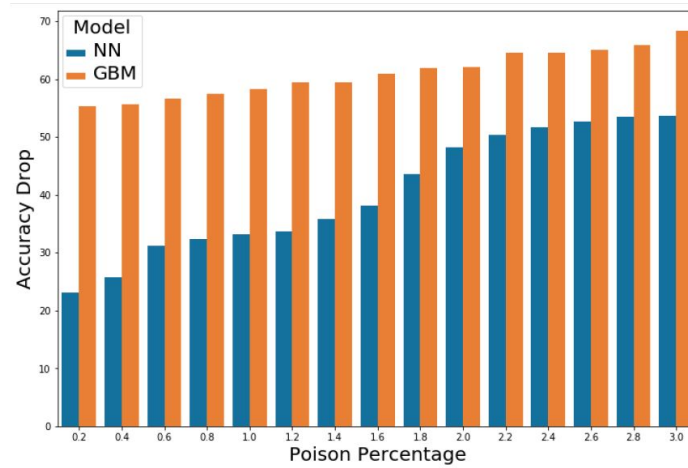
XAI to detect adversarial samples

- **Fidel et al. (2019)** use SHAP to extract the heatmaps of both genuine and adversarial samples
 - The heatmaps are used as **signatures** to detect adversarial samples



XAI to create adversarial samples

- **Kuppa et al. (2021)** propose a method to create adversarial sample with XAI:
 - Identify robust features that influence the class decision boundaries of the classifier
 - Perturbate robust features to create adversarial sample



Conclusion

Why Machine learning is vulnerable to adversarial attacks?

- Deep Neural network is data hungry
 - [Simon-Gabriel et al 2018 -Adversarial Vulnerability of Neural Networks Increases with Input Dimension]
- Nonlinearity of deep learning models
 - Overfitting of training data
 - Insufficient generalization ability to predict unknown data
- Adversarial vulnerability is a direct result of our models' sensitivity to well-generalizing features in the data:
 - Classifiers tend to use any available signal to do so, even those that look incomprehensible to humans
 - Different classifiers tend to find the same set of relevant features
 - that is why attacks can transfer across models!

Conclusion

- Defensive strategy:
 - Know your enemy
 - Attackers' goal and capability
 - Build machine learning model able to detect zero-days attacks
 - Reducing overfitting: defensive distillation, dropout
 - Design for attacks
 - Adversarial training
 - Feature robustness (e.g., XAI)

Classify this: is a cat or a dog?



References



- **Goodfellow, Ian J. et al.** "Deep Learning." *Nature* 521 (2015): 436-444. <https://www.nature.com/articles/nature14539>
- **Aggarwal, Charu C.** *Neural Networks and Deep Learning*. Cham: Springer, 2018.
- **GoodFellow, Ian et al.** "Deep Learning" <https://www.deeplearningbook.org/>
- **Rumelhart, D., Hinton, G. & Williams, R.** *Learning representations by back-propagating errors*. *Nature* 323, 533–536 (1986). <https://doi.org/10.1038/323533a0>
- **S. Albawi, T. A. Mohammed and S. Al-Zawi**, "Understanding of a convolutional neural network," *2017 International Conference on Engineering and Technology (ICET)*, 2017, pp. 1-6, doi: 10.1109/ICEngTechnol.2017.8308186
- **W. Wang, Y. Huang, Y. Wang and L. Wang**, "Generalized Autoencoder: A Neural Network Framework for Dimensionality Reduction," *2014 IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2014, pp. 496-503, doi: 10.1109/CVPRW.2014.79.
- **Goodfellow, Ian; Pouget-Abadie, Jean; Mirza, Mehdi; Xu, Bing; Warde-Farley, David; Ozair, Sherjil; Courville, Aaron; Bengio, Yoshua** (2014). *Generative Adversarial Nets* (PDF). Proceedings of the International Conference on Neural Information Processing Systems (NIPS 2014). pp. 2672–2680.
- **Ian J. Goodfellow, Jonathon Shlens, Christian Szegedy**: *Explaining and Harnessing Adversarial Examples*. ICLR (Poster) 2015
- **M. Melis, A. Demontis, B. Biggio, G. Brown, G. Fumera and F. Roli**, "Is Deep Learning Safe for Robot Vision? Adversarial Examples Against the iCub Humanoid," *2017 IEEE International Conference on Computer Vision Workshops (ICCVW)*, 2017, pp. 751-759, doi: 10.1109/ICCVW.2017.94
- **Huang Xiao, Battista Biggio, Blaine Nelson, Han Xiao, Claudia Eckert, and Fabio Roli**. 2015. *Support vector machines under adversarial label contamination*. *Neurocomput.* 160, C (July 2015), 53–62.
- **Ian J. Goodfellow, Jonathon Shlens, Christian Szegedy**, *Explaining and Harnessing Adversarial Examples*. ICLR (Poster) 2015
- **Kurakin, Alexey & Goodfellow, Ian & Bengio, Samy**. (2016). *Adversarial examples in the physical world*. Technical report, Google, Inc.
- **Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, Adrian Vladu**: *Towards Deep Learning Models Resistant to Adversarial Attacks*. ICLR (Poster) 2018

References



- **Rey Wiyatno, Anqi Xu:** *Maximal Jacobian-based Saliency Map Attack*. CoRR abs/1808.07945 (2018)
- **M. Al- Essa, G. Andresini, A. Appice, D. Malerba.,** *XAI to explore robustness of features in adversarial training for cybersecurity* , In ISMIS 2022. Lecture Notes in Computer Science(), vol 13515. Springer, Cham.
- **M. Al- Essa, G. Andresini, A. Appice, D. Malerba.,** *An XAI-based adversarial training approach for cyber-threat detection*. Accepted in *IEEE CyberSciTech/ PCom/ DASC/ CDBCom, 2022*