



Midterm 4

ON DETECTING ADVERSARIAL PERTURBATIONS

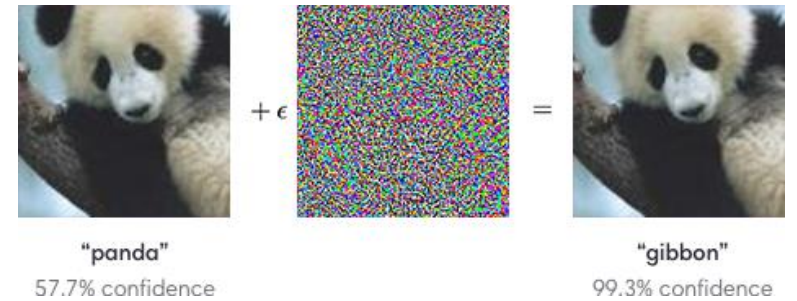
Intelligent Systems for Pattern Recognition
Master Degree in Computer Science, AI Curriculum
A.Y. 2020/2021

Introduction

Over the last years deep learning models obtained **impressive performances** in different complex tasks, above all:

- **Image classification** (Russakovsky et al., 2015; He et al., 2016) ;
- **Speech recognition** (Amodei et al., 2016);

Small and carefully directed perturbations are still the main concern for those tasks. The so called **adversarial examples** have been in the spotlight in the last years and many techniques have been developed recently for their implementation and to defend against them.



*“The vulnerability to adversarial inputs can be problematic and even prevent the application of deep learning methods in safety and security-critical applications, **the problem is particularly severe when human safety is involved.**”*

Many methods have been proposed for improving the robustness of networks against such attacks:

- **Augmenting the training data** with adversarial examples (Goodfellow et al., 2015);
- **Applying JPEG compression** to the input (Dziugaite et al., 2016);
- **Distilling a hardened network** from the original classifier network (Papernot et al., 2016);

The paper under examination propose and shows the training of a **binary detector to discriminate between non-adversarial and adversarial examples**, moreover it displays the generalization and transferability powers of some attacks and, finally, shows how an attack against both the classifier and detector can be developed and how to defend from it.

Adversarial attacks used

Fast method: The applied perturbation is the **direction** in image space which yields the **highest increase of the cost function**.

$$x^{adv} = x + \varepsilon \operatorname{sgn}(\nabla_x J_{cls}(x, y_{true}(x)))$$

Basic iterative method: Iterative version of the fast-method.

$$x_0^{adv} = x \quad x_{n+1}^{adv} = \operatorname{Clip}_x^\varepsilon \{x_n^{adv} + \alpha \operatorname{sgn}(\nabla_x J_{cls}(x_n^{adv}, y_{true}(x)))\} \quad L_\infty$$

$$x_0^{adv} = x \quad x_{n+1}^{adv} = \operatorname{Project}_x^\varepsilon \left\{ x_n^{adv} + \alpha \frac{\nabla_x J_{cls}(x_n^{adv}, y_{true}(x))}{\|\nabla_x J_{cls}(x_n^{adv}, y_{true}(x))\|} \right\} \quad L_2$$

- x input image;
- $y_{true}(x)$ one-hot encoding of the true class;
- $J_{cls}(x_n^{adv}, y_{true}(x))$ cost function of the classifier;
- $\sigma \in [0, 1]$ hyperparameter that works as a tradeoff;
- $J_{det}(x, 1)$ cost function of the detector.

DeepFool method: It iteratively perturbs an image x , at each step the classifier is linearized around the n -th x and the closest class boundary is determined. The **attack stops when the actual class given by the classifier changes**. The paper uses both the L_∞ and L_2 norm.

*“In the worst case, an adversary might not only have access to the classification network and its gradient **but also to the adversary detector and its gradient**. In this case, the adversary might potentially generate inputs to the network that fool both the classifier (i.e., **get classified wrongly**) and fool the detector (i.e., **look innocuous**).”*

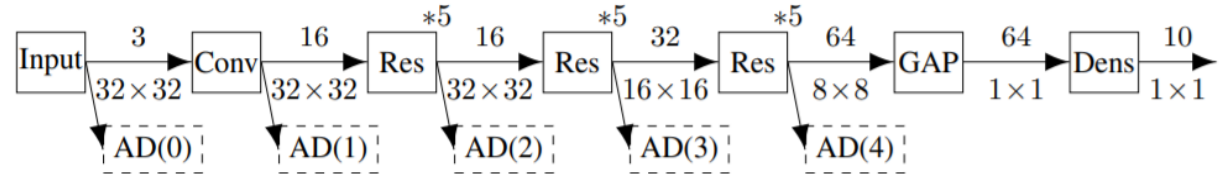
Dynamic attack:

$$x_0^{adv} = x \quad x_{n+1}^{adv} = \operatorname{Clip}_x^\varepsilon \left\{ x_n^{adv} + \alpha \left((1 - \sigma) \operatorname{sgn}(\nabla_x J_{cls}(x_n^{adv}, y_{true}(x))) + \sigma \operatorname{sgn}(\nabla_x J_{det}(x_n^{adv}, 1)) \right) \right\}$$

Training and test with Cifar-10

Classifier:

- **Epochs:** 100;
- **Momentum:** 0.9 (starting);
- **Learning rate:** 0.1, 0.01 (after 41 epochs), 0.001 (after 61 epochs);
- Uses **stochastic gradient descent**;
- Network **performance are evaluated after each epoch** on the validation data.



Accuracy on test set: 91,3%

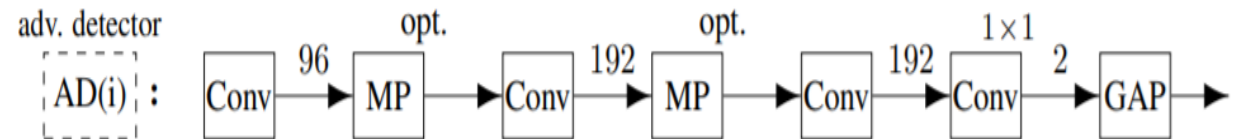
First we **train the classifier** on the regular dataset;

Then we **augment the initial dataset by generating an adversarial example for each image**, doubling the size of the dataset;

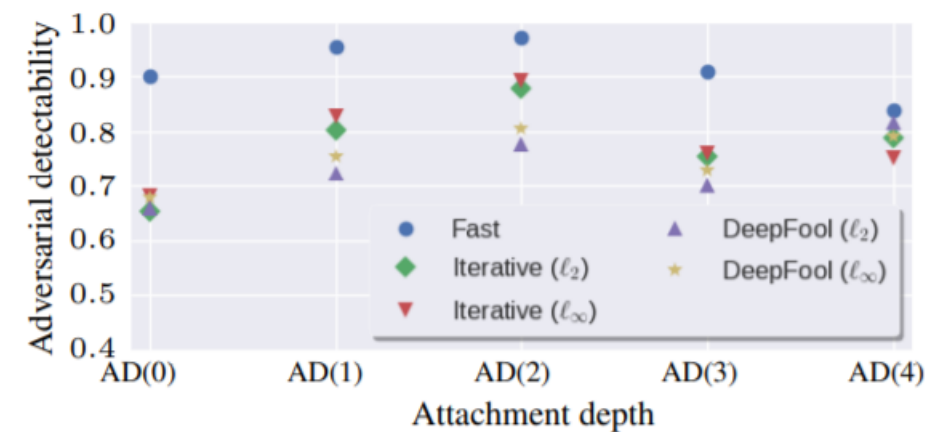
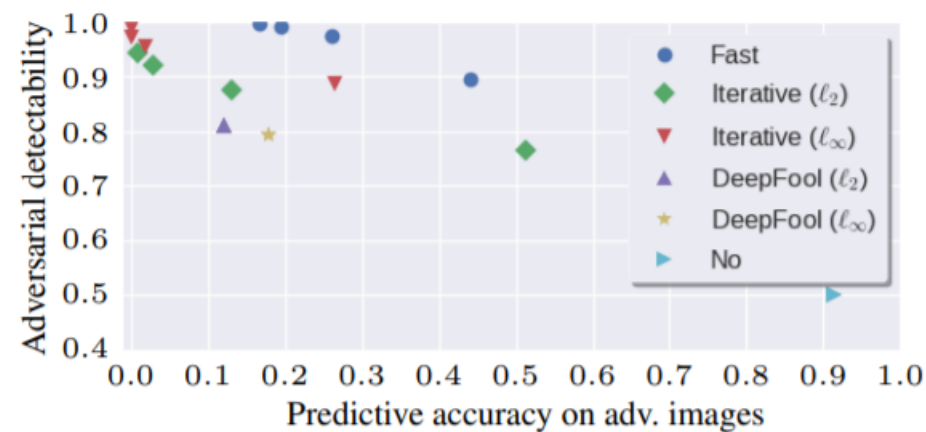
Freeze the weights of the classifier and **train the detector**.

Detector:

- **Epochs:** 20;
- **Optimizer:** Adam;
- **Learning rate:** 0.0001;
- $\beta_1 = 0.99$, $\beta_2 = 0.999$;
- Produces a value $p \in [0, 1]$;
- The detector was **attached to position AD(2)** for **each adversary attack** except for DeepFool where the detector was **attached on AD(4)**, this will be explained in the next slide.

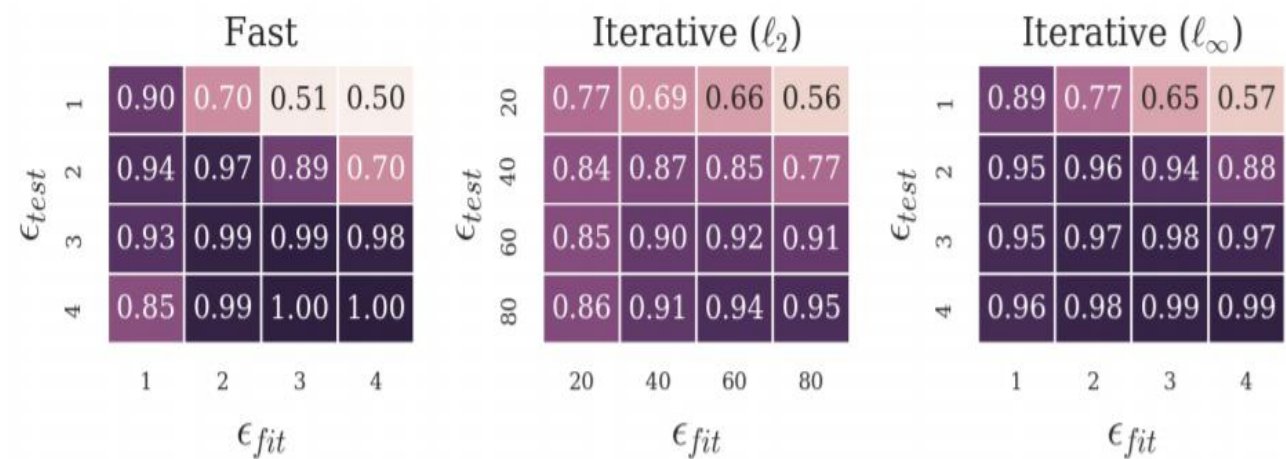


Results on Cifar-10

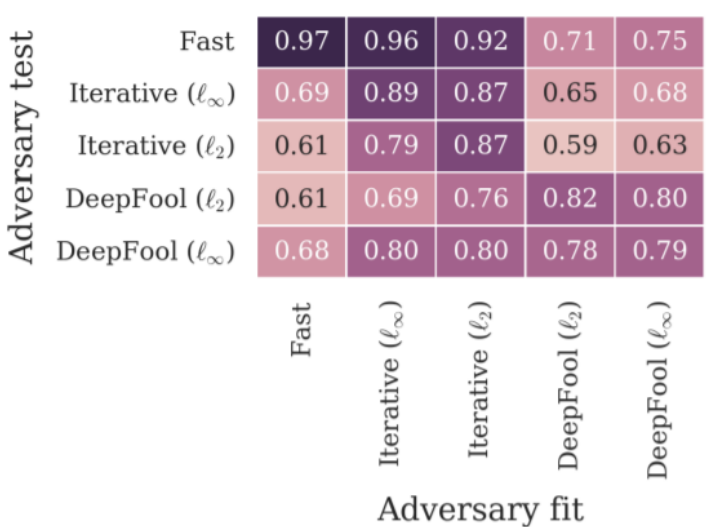


The “Fast” adversary can be considered as a weak adversary, the DeepFool based methods as relatively strong adversaries, and the “Iterative” method being somewhere in-between.

For the “Fast” and “Iterative” adversaries AD(2) is the best, meanwhile AD(4) works best for DeepFool.



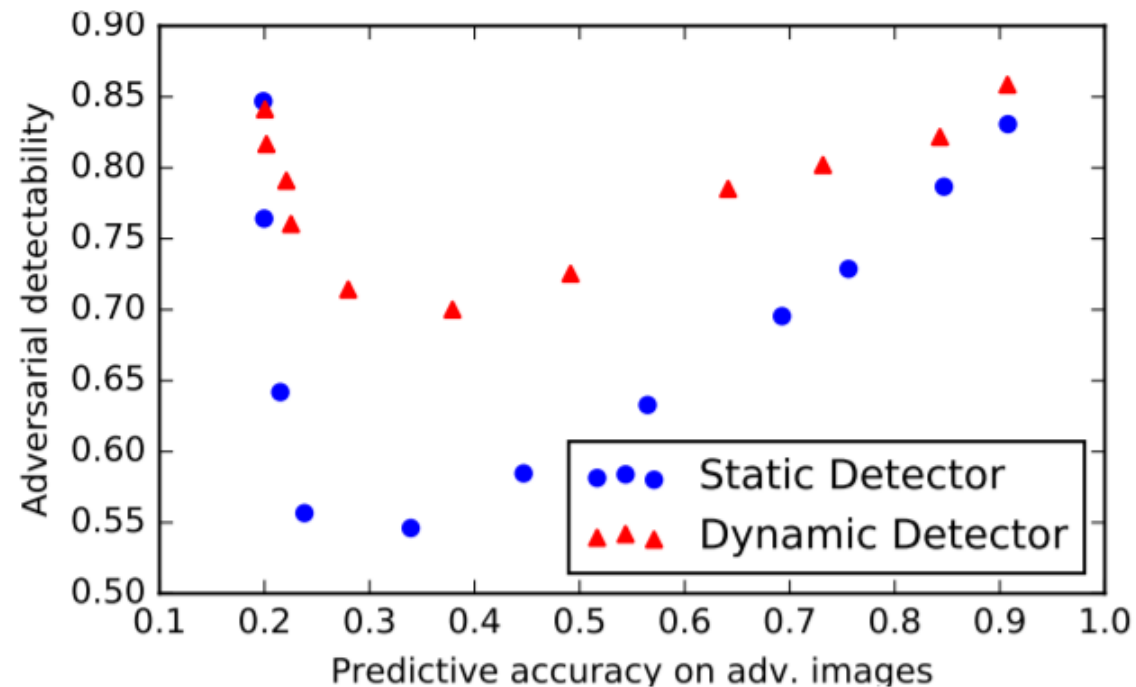
Transferability on same adversary with different ϵ



Transferability on different adversaries

Dynamic detector on Cifar-10

Dynamic adversary training: instead of precomputing a dataset of adversarial examples, we compute the adversarial examples on-the-fly for each mini-batch and let the adversary modify each data point with probability 0.5. Note that a dynamic adversary will modify a data point differently every time it encounters the data point since it depends on the detector's gradient and the detector changes over time. By training the detector in this way, we implicitly train it to resist dynamic adversaries for various values of σ .

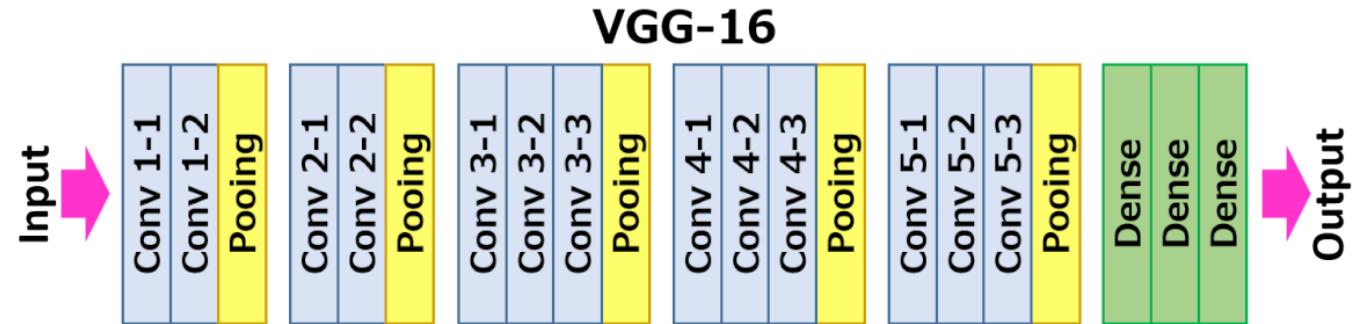


As we can see from the image static detectors performed worse than dynamic ones (as expected) on dynamic adversarial attacks, the parameter σ has been chosen as $\sigma \in \{0.0, 0.1, \dots, 1.0\}$, with smaller values of σ corresponding to lower predictive accuracy, i.e., being further on the left.

10-class ImageNet

Classifier:

- Pretrained VGG16 (Simonyan & Zisserman, 2015);
- Before the SoftMax another layer is placed so that we can restrict the classification to only ten classes.

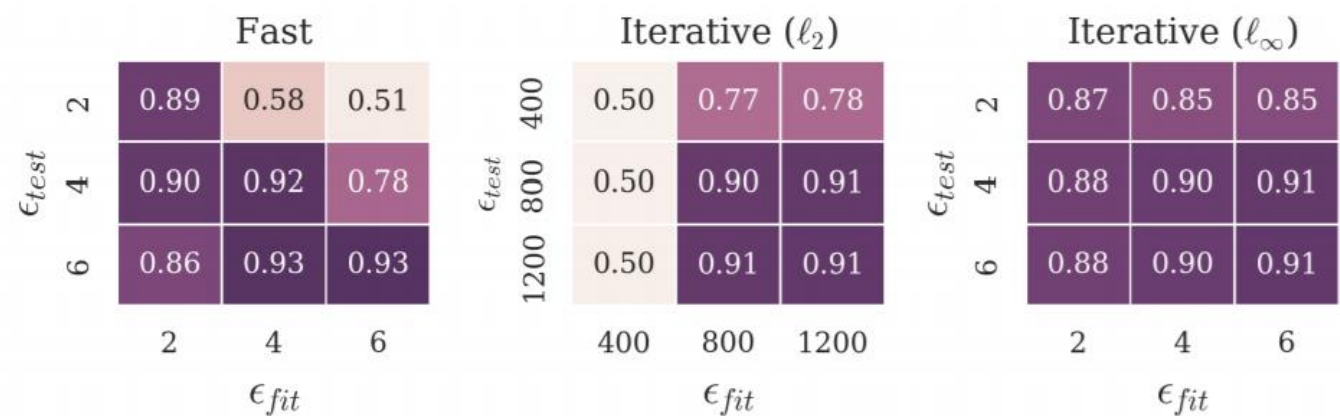
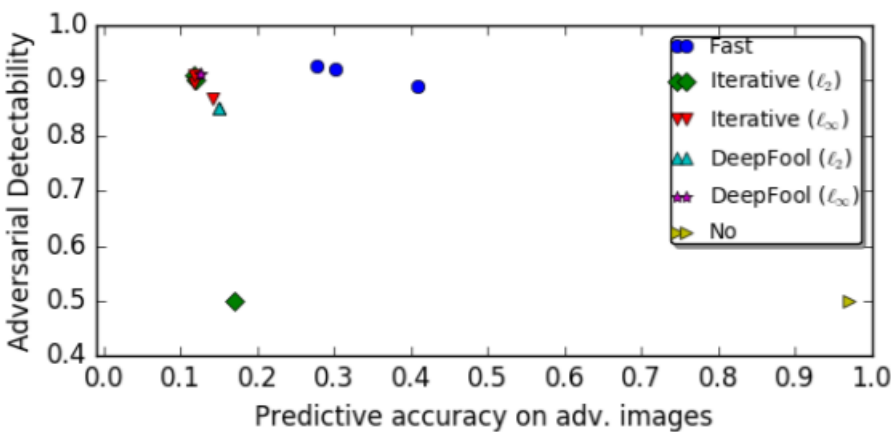


Detector:

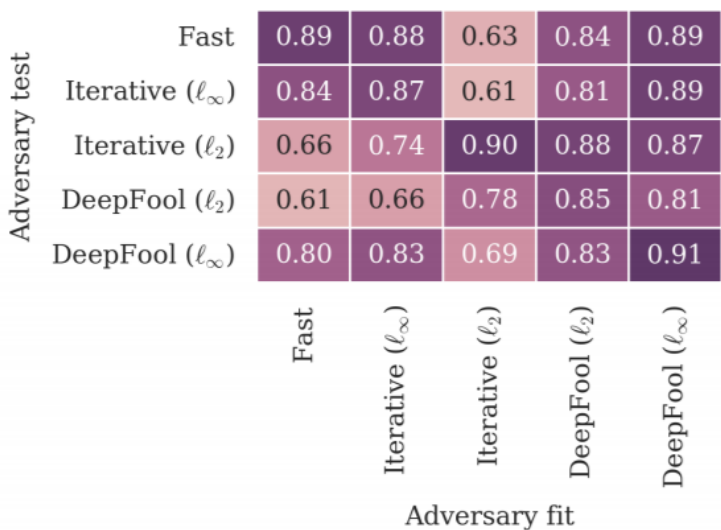
- **Epochs:** 500;
- **Optimizer:** Adam;
- **Learning rate:** 0.0001;
- $\beta_1 = 0.99$, $\beta_2 = 0.999$;
- Produces a value $p \in [0, 1]$;
- The detector was attached after the fourth max pooling layer;
- The network consists of a **sequence of five 3x3 convolutions** with 196 feature maps each using **batch-normalization** and **rectified linear units**, followed by a 1x1 convolution which maps onto the 10 classes, **global-average pooling**, and a **softmax layer**. Another 2x2 is inserted after the first convolution

Results on 10-class ImageNet

The image shows that detectability is 85% percent or more with the exception of the “Iterative” ℓ_2 -based adversary with $\epsilon = 400$. For this adversary, the detector only reaches chance level.



Transferability on same adversary with different ϵ , the results are similar to the ones obtained on the cifar-10 dataset



Transferability on different adversaries

Conclusions

- The results shown on the paper are somewhat complementary to Moosavi-Dezfooli et al. (2016a): while they show that universal, image-agnostic perturbations exist, the paper has shown that image-dependent perturbations are sufficiently regular to be detectable.
- **Transferability** is not symmetric and typically **works best between similar adversaries and from stronger to weaker adversary**;
- For a **static detector**, there might be areas which are **adversarial to both classifier and detector**; however, this will be a (small) subset of the areas which are adversarial to the classifier alone. However, a **dynamic detector is considerably harder to fool**: on the one hand, it might further **reduce the number of areas which are both adversarial to classifier and detector**.

*The paper has shown that **adversarial examples can be detected surprisingly well using a detector subnetwork** attached to the main classification network, **this does not directly allow classifying adversarial examples correctly**, but **it allows mitigating adversarial attacks against machine learning systems** by resorting to **fallback solutions**. Additional future work will be developing **stronger adversaries that are harder to detect** by adding effective randomization which would make selection of adversarial perturbations less regular. Finally, **developing methods for training detectors** such that they can **detect many different kinds of attacks reliably at the same time** would be essential for safety and security-related applications.*

Thanks for the attention