

A tutorial on Deep Learning for Privacy in Multimedia

Part 2: Adversarial images

Ali Shahin Shamsabadi

a.shahinshamsabadi@qmul.ac.uk







What is an adversarial image?





- Unnoticeable
 - Humans do not perceive any visual distortions and keep aspect ratio



Adversarial Image[CVPR-W'17]



Adversarial Image[CVPR'20]







CVPR-W'17: Hosseini and Poovendran, "Semantic adversarial examples". CVPR'20: Shahin Shamsabadi et al, "ColorFool: Semantic adversarial colorization".

An adversarial image should be

- Unnoticeable
 - Humans do not perceive any visual distortions and keep aspect ratio
- Transferable
 - The adversarial image can mislead unseen classifiers (i.e. do not overfit to one classifier)





An adversarial image should be

- Unnoticeable
 - Humans do not perceive any visual distortions and keep aspect ratio
- Transferable
 - The adversarial image can mislead unseen classifiers (i.e. do not overfit to one classifier)
- Robust
 - The perturbation is not removed by defence frameworks[NDSS'18]
 - Quantization, Median smoothing and JPEG compression



An adversarial image should be

- Unnoticeable
 - Humans do not perceive any visual distortions and keep aspect ratio
- Transferable
 - The adversarial image can mislead unseen classifiers (i.e. do not overfit to one classifier)
- Robust
 - The perturbation is not removed by defense frameworks[NDSS'18]
 - Quantization, Median smoothing and JPEG compression

	Knowledge	Example	
White-box	Everything	Open source classifiers	Best-case scenario for the attacker
Black-box	Predictions	Public APIs	-
No box	Nothing	Classifiers on social media	Worst-case scenario for the attacker



- Untargeted
 - Mislead with any class that is different from the class of the original image

- Targeted
 - Mislead to predict a specific target class









Generating an adversarial image against $C(\cdot)$



Fixed parameters, Adjust input

$$\dot{X} = \operatorname{argmax}_{X} J(C(\dot{X}), y)$$

$$\dot{X}$$
Adversarial image





Generating an adversarial image against $C(\cdot)$



$$\dot{X} = \underset{\dot{X}}{\operatorname{argmax}} J(C(\dot{X}), y) \text{ s.t. } \begin{cases} \dot{X} \text{ is still an image} \\ \operatorname{Perceptually similar} \end{cases}$$



CIS centre for intelligent sensing

Generating an adversarial image against $C(\cdot)$

• Parameters are fixed, maximising the gradients of the loss function wrt input



$$\dot{X} = \operatorname*{argmax}_{\dot{X}} J(C(\dot{X}), y) \text{ s.t.} \begin{cases} \dot{X} \text{ is still an image} \\ \operatorname{Perceptually similar} \end{cases}$$



So how we perturb images?

- Possible space of allowed perturbations
 - Norm-bounded perturbations
 - Content-based perturbations



Norm-bounded perturbations

- Minimize l_p -norm
 - Maximum change for each pixel, l_{∞}
 - Fast Gradient Sign Method (FGSM)[ICLR'14]
 - Basic Iterative Method (BIM)[ICLR'17]
 - Robust and private BIM (RP-BIM)[ТММ'20]
 - Maximum energy change, l_2
 - DeepFool (l_2) [CVPR'16]
 - Carlini-Wagner (CW)[S&P'17]
 - Maximum number of perturbed pixels, l_1
 - SparseFool[CVPR'19]
 - JSMA[EuroS&P'16]
- ICLR'14: Goodfellow et al, "Explaining and harnessing adversarial examples".
- ICLR'17: Kurakin et al, "Adversarial examples in the physical world".
- TMM'20: Sanchez-Matilla et al, "Exploiting vulnerabilities of deep neural networks for privacy protection".
- CVPR'16: Moosavi-Dezfool et al, "DeepFool: A simple and accurate method to fool deep neural networks ".
- S&P'17: Carlini and Wagner, "Towards evaluating the robustness of neural networks".
- CVPR'19: Modas et al, "SparseFool: a few pixels make a big difference".
- EuroS&P'16: Papernot et al. "The limitations of deep learning in adversarial settings".













Fast Gradient Sign Method, untargeted



ICLR'14: Goodfellow et al, "Explaining and harnessing adversarial examples".

Fast Gradient Sign Method, targeted



ICLR'14: Goodfellow et al, "Explaining and harnessing adversarial examples".

Improving misleading of Fast Gradient Sign Method

- Basic Iterative Method: Iteratively generate the adversarial perturbation
- Aggregates N perturbations

$$\dot{X}_0 = X$$
$$\dot{X} = \operatorname{Clip}_{\epsilon}(X + \sum_{n} \operatorname{sign}\left(\nabla_{\dot{X}_{n-1}} J(C(\dot{X}_{n-1}), y)\right))$$



Improving robustness and transferability of FGSM/BIM



TMM'20: Sanchez-Matilla et al, "Exploiting vulnerabilities of deep neural networks for privacy protection".

In summary, norm-bounded perturbations are

- Unnoticeable to human eyes
 - Small perturbations
- Not transferable to unseen classifiers
 - Overfitted to a specific classifier
 - Small magnitudes
- Not robust to defenses
 - High-frequency



In summary, norm-bounded perturbations are

- Unnoticeable to human eyes
 - Small perturbations
- Not transferable to unseen classifiers
 - Overfitted to a specific classifier
 - Small magnitudes
- Not robust to defenses
 - High-frequency

Neglect image content



How image contents can help adversarial perturbations?

- Introduce a larger range of perturbations based on image content
 - Colour
 - SemanticAdv[CVPR-W'17]
 - ColorFool[CVPR'20]
 - ACE[BMVC'20]
- Structure and details
 - EdgeFool[ICASSP'20]
 - FilterFool[Arxiv'20]

CVPR-W'17: Hosseini and Poovendran, "Semantic adversarial examples".

CVPR'20: Shahin Shamsabadi et al, "ColorFool: Semantic adversarial colorization".

BMVC'20: Zhao et al, "Adversarial Color Enhancement: Generating Unrestricted Adversarial Images by Optimizing a Color Filter".

ICASSP'20: Shahin Shamsabadi et al, "EdgeFool: An adversarial image enhancement filter ".

Arxiv'20: Shahin Shamsabadi et al, "Semantically Adversarial Learnable Filters ".



Large colour perturbations

- Untargeted adversarial colour changes
 - HSV colour space
 - Shifting hue and saturation
- Low-frequency colour perturbations
 - Transferable
 - Robust
 - Black-box
 - Unnatural-looking adversarial images

SemanticAdv

Original









CVPR-W'17: Hosseini and Poovendran, "Semantic adversarial examples".

Why adversarial colours may look unnatural?

- Images contain semantic regions
 - A set of identifiable pixels such as car or person
- Sensitive regions
 - Colour changes are noticeable
 - Person, Vegetation, Water and Sky
- Non-sensitive regions
 - Occur in a wide range of colours
 - Curtain, Wall and ...













We also need a better perceptually motivated colour space

- Lab colour space
- Perceptually uniform
- Represent colour separately from lightness
 - Lightness
 - L channel: black (0) to white (100)
 - Colour information
 - a channel: green (-128) to red (+128)
 - b channel: blue (-128) to yellow (+128)





Semantic and perceptually motivated colour perturbations

- Perform in Lab colour space
- Modify de-correlated a and b colour channels
 - Semantic regions
 - Human perception
- Up to 1000 trials



co_cop_QuantifyingColors.html

Semantic and perceptually motivated colour perturbations

- Perform in Lab colour space
- Modify a and b colour channels
 - Semantic regions
 - Human perception
- Up to 1000 trials

Original









ColorFool-ed

ColorFool

- Transferable
 - Randomness
 - Black-box
- Robust
 - Low frequency
- Natural-looking
 - Semantic regions
 - Human perception
- Maintain aspect ratio





Now let's see how we can implement ColorFool-ed images

Original





https://github.com/smartcameras/ColorFool

What about other semantic attributes of images?





ICASSP'20: Shahin Shamsabadi et al, "EdgeFool: An adversarial image enhancement filter ".

Can perturbations perform enhancement, not distortion?

- Yes
- Exploiting traditional image processing filters



ICASSP'20: Shahin Shamsabadi et al, "EdgeFool: An adversarial image enhancement filter ".



ICASSP'20: Shahin Shamsabadi et al, "EdgeFool: An adversarial image enhancement filter".



Enhance image detail in the Lab space to not distort the color

 $\dot{\mathbf{X}} = [f(\mathbf{X}_d^L, \mathbf{X}_s^L), \mathbf{X}^a, \mathbf{X}^b]$

ICASSP'20: Shahin Shamsabadi et al, "EdgeFool: An adversarial image enhancement filter ".



Guide the enhancement in an adversary manner via minimizing

$$(L_a) = (\dot{z}_y) - \max\{\dot{z}_i, i \neq y\}$$

Adversarial Loss Function Score of belonging adversarial image to the same class as the original image

ζ

Classifier

ICASSP'20: Shahin Shamsabadi et al, "EdgeFool: An adversarial image enhancement filter ".

Classifier



Residual learning: if you look after other enhancements

- Enhanced image outputted by the filter: X_e
- The filter residual: $\delta_e = X_e X$
- Learnable adversarial residual: δ



Penalizes pixel-wise large differences Guides the residual towards adversarial



Arxiv'20: Shahin Shamsabadi et al, "Semantically Adversarial Learnable Filters".

Original images

FilterFool-ed



Detail enhanced





Gamma corrected





Log transformed



Others traditional image processing filters

- Adversary learn parameters of filters to output adversarial images
- Differentiable approximation of a colour filter

$$\dot{X}_{\chi} = \sum_{i=1}^{\left[\frac{X_{\chi}}{S}\right]} \theta_{i} + \frac{X_{\chi}(\text{mod}s)}{s} \theta_{\left[\frac{X_{\chi}}{S}\right]}$$





Now let's look at class labels



? "Television" "Typewriter keyboard"

? "collie"



Now let's look at class labels



Adversarial:

Original:

"Television"

"Typewriter keyboard"

"collie"



Arxiv'20: Shahin Shamsabadi et al, "Semantically Adversarial Learnable Filters".

Semantic relationships between class labels

- Cluster D = 1000 ImageNet classes to S = 11 semantic classes
 - Dogs (130); Other mammals (88); Birds (59); Reptiles, fish, amphibians (60);
 Invertebrates (61); Food, plants, fungi (63); Devices (172); Structures, furnishing (90);
 Clothes, covering (92); Implements, containers, misc. objects (117); and Vehicles (68)
- Define mapping matrix $W \in \{0,1\}^{D \times S}$

Semantically different classes than the original class

$$L_{S-Adv}(X, \dot{X}) = < \operatorname{ReLU}(\dot{z}), w_{s} > -\max(\dot{z} \odot \widehat{w}_{s})$$
Semantically similar classes

to the original class

 $<\cdot,\cdot>$: dot product \odot : Hadamard product

Semantic relationships between class labels



Dogs (130);
Other mammals (88);
Birds (59);
Reptiles, fish, amphibians (60);
Invertebrates (61);
Food, plants, fungi (63);
Devices (172);
Structures, furnishing (90);
Clothes, covering (92);
Implements, containers, misc. objects (117);
Vehicles (68).

$$\mathcal{L}_{S-Adv}(\boldsymbol{I}, \boldsymbol{\dot{I}}) = < \text{ReLU}(\boldsymbol{\dot{z}}), \boldsymbol{w}_{s} > -\max\{\boldsymbol{\dot{z}} \odot \boldsymbol{\widehat{w}}_{s}\}$$

Experimental settings

- Dataset
 - ImageNet: 3K images of 1k objects
- Classifiers under attacks
 - ResNet50, ResNet18, AlexNet
- Visualization
 - Perturbations
 - Adversarial images
 - Top5 predictions and classifier attention
- Success rate and transferability

successful adversarial images

CIS centre for intelligent sensing # total images



Adversarial perturbations

Original





SparseFool



SemanticAdv





EdgeFool

FilterFool



Adversarial images



Basic Iterative Method



ColorFool



SparseFool



EdgeFool

SemanticAdv





Top-5 predictions

Original



Basic Iterative Method



SparseFool



SemanticAdv



ColorFool



EdgeFool collie

Shetland sheepdog Australian terrier German shepherd Border collie

FilterFool



sunscreen koala sunglasses

Classifier attention

Original



Basic Iterative Method



ColorFool



SparseFool



SemanticAdv



FilterFool





ResNet50 (seen) (•) ResNet18 (unseen) (•) AlexNet (unseen) (•)

BIM: Basic Iterative Method, SF: SparseFool, SA: SemanticAdv,EF: EdgeFool, FF: FilterFool, LT: Log transformation,LD, ND: Linear and Nonlinear Detail enhancement,G: Gamma correction

- Images are special data not like tabular data
- Improves unnoticeability, transferability and robustness
 - Large content-based perturbations
 - Structure, detail, colours and objects

• Protect privacy in photo sharing social media





Still much room for improvement!

- Define mathematically the possible set of adversarial images
 - Model human eyes
- Why classifiers are vulnerable to semantic changes?
 - Colour shifting
 - Detail enhancement
- Certified adversarial perturbations
- Make semantic changes outside of digital domains in physical world



