



# Advanced Computer Vision

Практический курс

Савельева Юлия Олеговна, 3 семестр, 03.12.2020

# Adversarial Examples

Взять уже обученную на Stanford Online Products (from scratch, не pretrained на ImageNet) нейронную сеть на задачу классификации.  
Реализовать FGSM или T-FGSM, **НО не** I-FGSM.

## Fast gradient sign method (FGSM)

This method computes an adversarial image by adding a pixel-wide perturbation of magnitude in the direction of the gradient. This perturbation is computed with a single step, thus is very efficient in terms of computation time:

$$x^{adv} = x + \varepsilon \cdot \text{sign}(\nabla_x J(x, y_{true})),$$

where

$x$  is the input (clean) image,

$x^{adv}$  is the perturbed adversarial image,

$J$  is the classification loss function,

$y_{true}$  is true label for the input  $x$ .

<https://medium.com/onfido-tech/adversarial-attacks-and-defences-for-convolutional-neural-networks-66915ece52e7>

# Adversarial Examples

## Targeted fast gradient sign method (T-FGSM)

Similarly to the FGSM, in this method a gradient step is computed, but in this case in the direction of the negative gradient with respect to the target class:

$$x^{adv} = x - \epsilon \cdot \text{sign}(\nabla_x J(x, y_{target})),$$

where

$y_{target}$  is the target label for the adversarial attack.

## Iterative fast gradient sign method (I-FGSM)

The iterative methods take  $T$  gradient steps of magiture  $\alpha = \epsilon / T$  instead of a single step  $t$ :

$$x_0^{adv} = x, \quad x_{t+1}^{adv} = x_t^{adv} + \alpha \cdot \text{sign}(\nabla_x J(x_t^{adv}, y)).$$

Both one-shot methods (FGSM and T-FGSM) have lower success rates when compared to the iterative methods (I-FGSM) in white box attacks, however when it comes to black box attacks the basic single-shot methods turn out to be more effective. The most likely explanation for this is that the iterative methods tend to overfit to a particular model.

# Adversarial Examples

Для картинок из списка  
вычислить adversarial examples.  
Найти такой  $\epsilon$ , чтобы нейронная  
сеть справлялась с  
классификаций хуже всего, но  
при этом на картинке не  
появлялось видимого шума.  
Изменить картинки и прислать  
на проверку.

bicycle\_final/111265348817\_0.JPG  
bicycle\_final/111265348817\_1.JPG  
bicycle\_final/111265348817\_2.JPG  
bicycle\_final/111265348817\_3.JPG  
bicycle\_final/111469262153\_0.JPG  
bicycle\_final/111588452395\_0.JPG  
bicycle\_final/111588452395\_1.JPG  
bicycle\_final/111612717975\_0.JPG  
bicycle\_final/111612717975\_1.JPG  
bicycle\_final/111612717975\_2.JPG  
bicycle\_final/111612717975\_3.JPG  
bicycle\_final/111612717975\_5.JPG  
bicycle\_final/111620439270\_1.JPG  
bicycle\_final/111620439270\_3.JPG  
bicycle\_final/111645993618\_0.JPG  
bicycle\_final/111645993618\_1.JPG  
bicycle\_final/111645993618\_2.JPG  
bicycle\_final/111645993618\_3.JPG  
bicycle\_final/111661598505\_0.JPG  
bicycle\_final/111661598505\_2.JPG  
bicycle\_final/111661598505\_3.JPG  
bicycle\_final/111661603079\_0.JPG  
bicycle\_final/111661603079\_2.JPG  
bicycle\_final/111661603079\_3.JPG  
bicycle\_final/111687027218\_0.JPG

bicycle\_final/111687027218\_1.JPG  
bicycle\_final/111687027218\_10.JPG  
bicycle\_final/111687027218\_11.JPG  
bicycle\_final/111687027218\_2.JPG  
bicycle\_final/111687027218\_3.JPG  
bicycle\_final/111687027218\_4.JPG  
bicycle\_final/111687027218\_6.JPG  
bicycle\_final/111687027218\_7.JPG  
bicycle\_final/111687027218\_8.JPG  
bicycle\_final/111687027218\_9.JPG  
bicycle\_final/111688507119\_1.JPG  
bicycle\_final/111688507119\_2.JPG  
bicycle\_final/111702172753\_0.JPG  
bicycle\_final/111702172753\_1.JPG  
bicycle\_final/111702172753\_2.JPG  
bicycle\_final/111708813226\_0.JPG  
bicycle\_final/111708813226\_1.JPG  
bicycle\_final/111720342775\_0.JPG  
bicycle\_final/111720342775\_1.JPG  
bicycle\_final/111720342775\_10.JPG  
bicycle\_final/111720342775\_11.JPG  
bicycle\_final/111720342775\_3.JPG  
bicycle\_final/111720342775\_4.JPG  
bicycle\_final/111720342775\_5.JPG  
bicycle\_final/111720342775\_8.JPG

Дедлайн 10.12.2020 03:00