

Adversarial examples

Гельван Кирилл

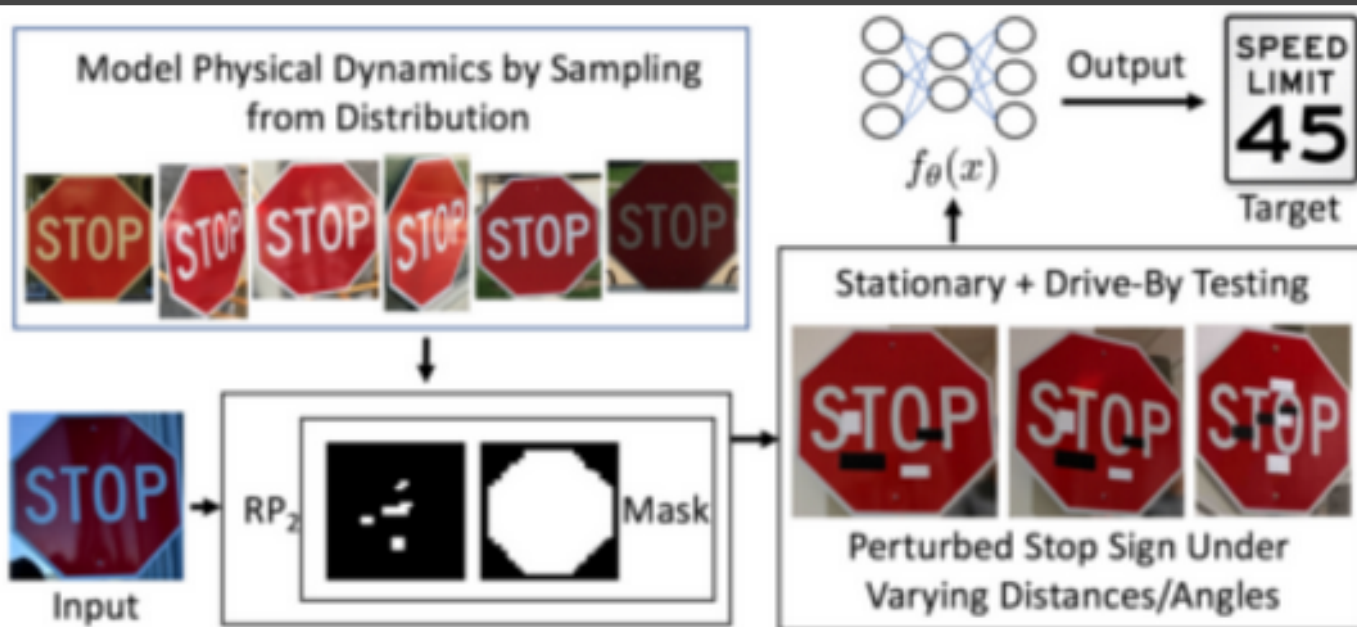
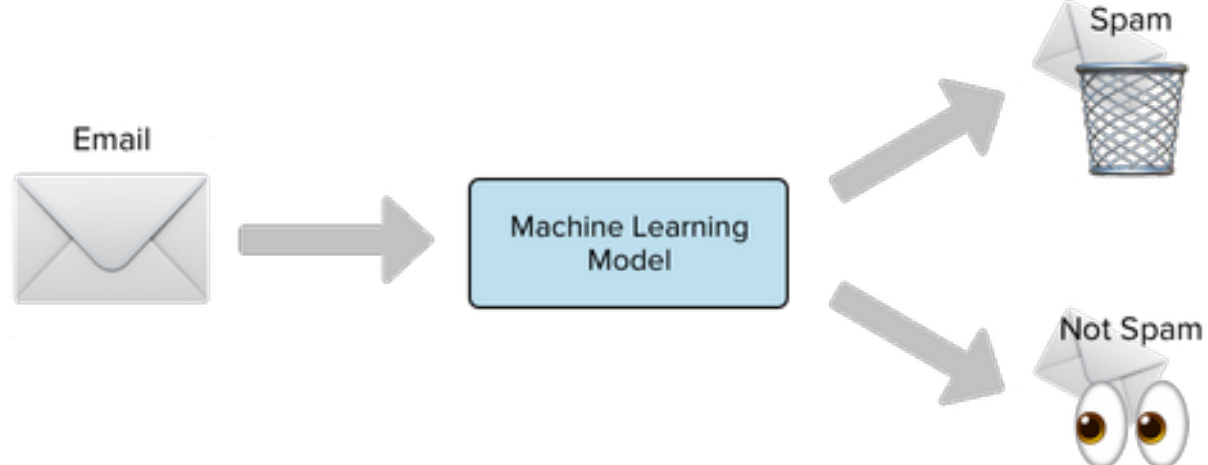
НИУ ВШЭ

15.11.19

Что сейчас будет?

- I. Что такое adversarial examples?
- I. Как создавать adversarial examples?
- I. Как защищаться от adversarial атак?

I. Что такое adversarial examples?



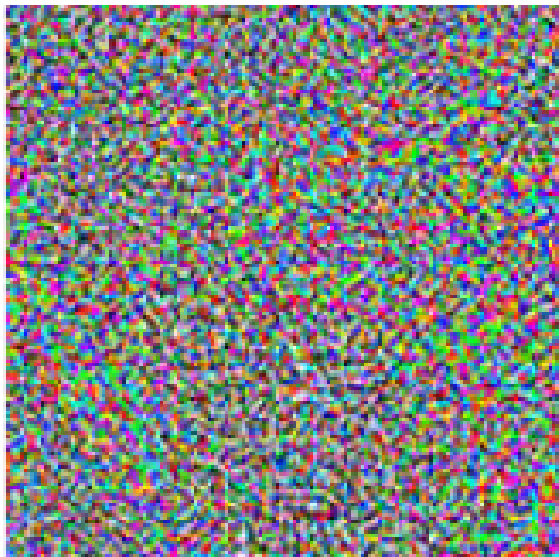
Оригинал



"panda"

57.7% confidence

+ ϵ



=



"gibbon"

99.3% confidence



Наложение на картинку, может заставить классификатор определить панду в категорию «гibbon»

Обозначение	Значение
x	оригинальные данные
l	номер класса ($1 \dots m$)
x'	adversarial example
l'	номер класса (если adv. ex. направ.)
$f(\cdot)$	модель ($f \in F : \mathbb{R}^n \rightarrow l$)
θ	параметры модели f
$J_f(\cdot, \cdot)$	функция потерь для модели f
η	разница между оригиналом и adv. ex. ($\eta = x' - x$)
$\ \cdot\ _p$	l_p норма
∇_t	градиент по t

Постановка задачи

Дано:

Модель f , оригинал x

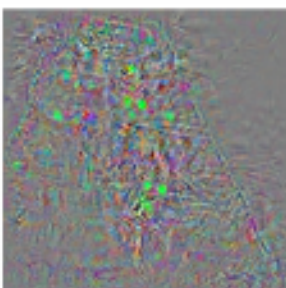
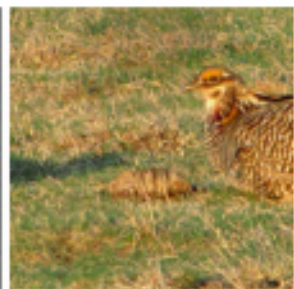
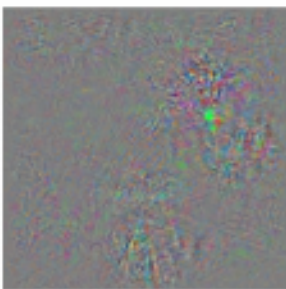
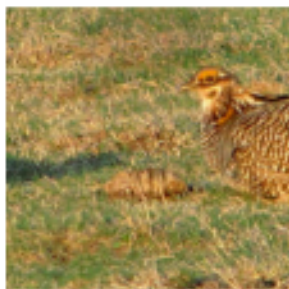
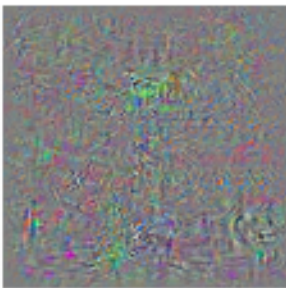
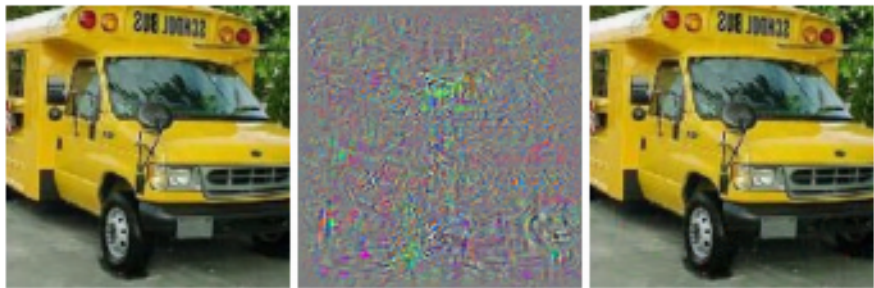
Задача:

$$\begin{aligned} \min_{x'} \quad & \|x' - x\| \\ \text{s.t.} \quad & f(x') = l', \\ & f(x) = l, \\ & l \neq l', \\ & x' \in [0, 1], \end{aligned}$$

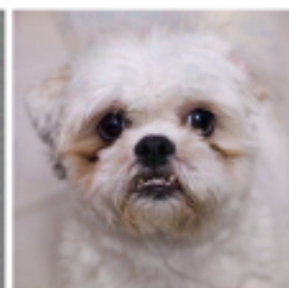
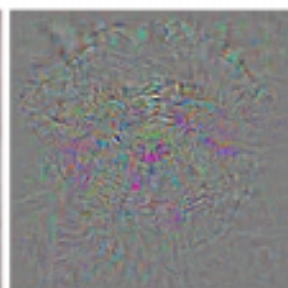
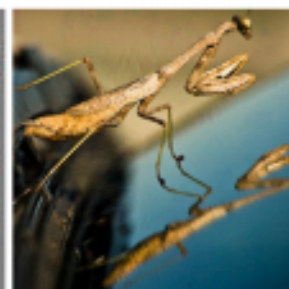
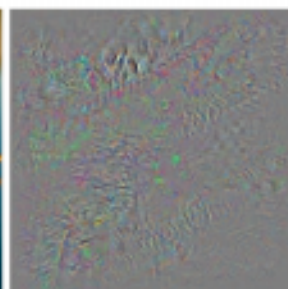
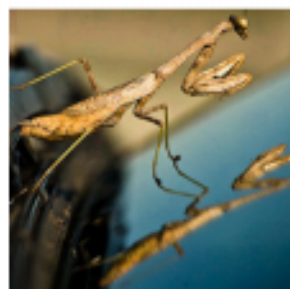
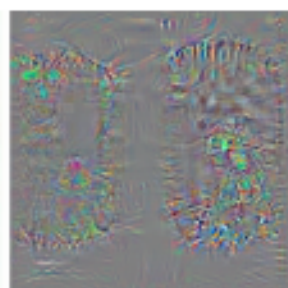
$$\eta = x' - x$$

Perturbation ~ разница,
отклонение

Обозначение	Значение
x	оригинальные данные
l	номер класса ($1 \dots m$)
x'	adversarial example
l'	номер класса (если adv. ex. направ.)
$f(\cdot)$	модель ($f \in F : \mathbb{R}^n \rightarrow l$)
θ	параметры модели f
$J_f(\cdot, \cdot)$	функция потерь для модели f
η	разница между оригиналом и adv. ex. ($\eta = x' - x$)
$\ \cdot\ _p$	l_p норма
∇_t	градиент по t



$$x + \epsilon x' = \text{страус}$$



$$x + \epsilon x' = \text{страус}$$

Threat model

Adversarial
Falsification

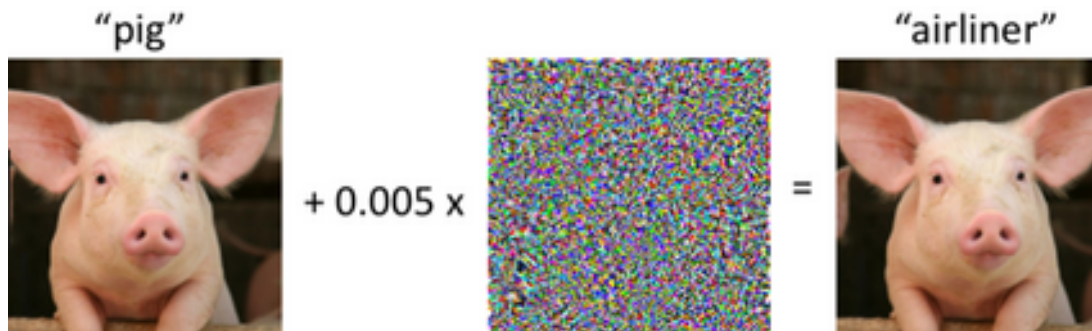
Adversarial
Specificity

Adversary's
Knowledge

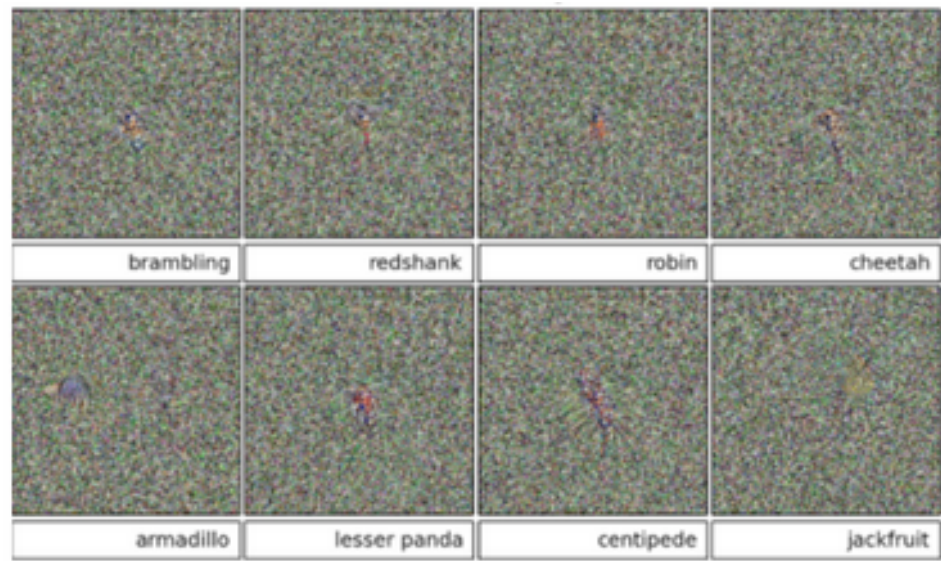
Attack Frequency

Adversarial Falsification

False Negative



False Positive



Уверенность $\geq 99\%$

Adversary's Knowledge

White-box



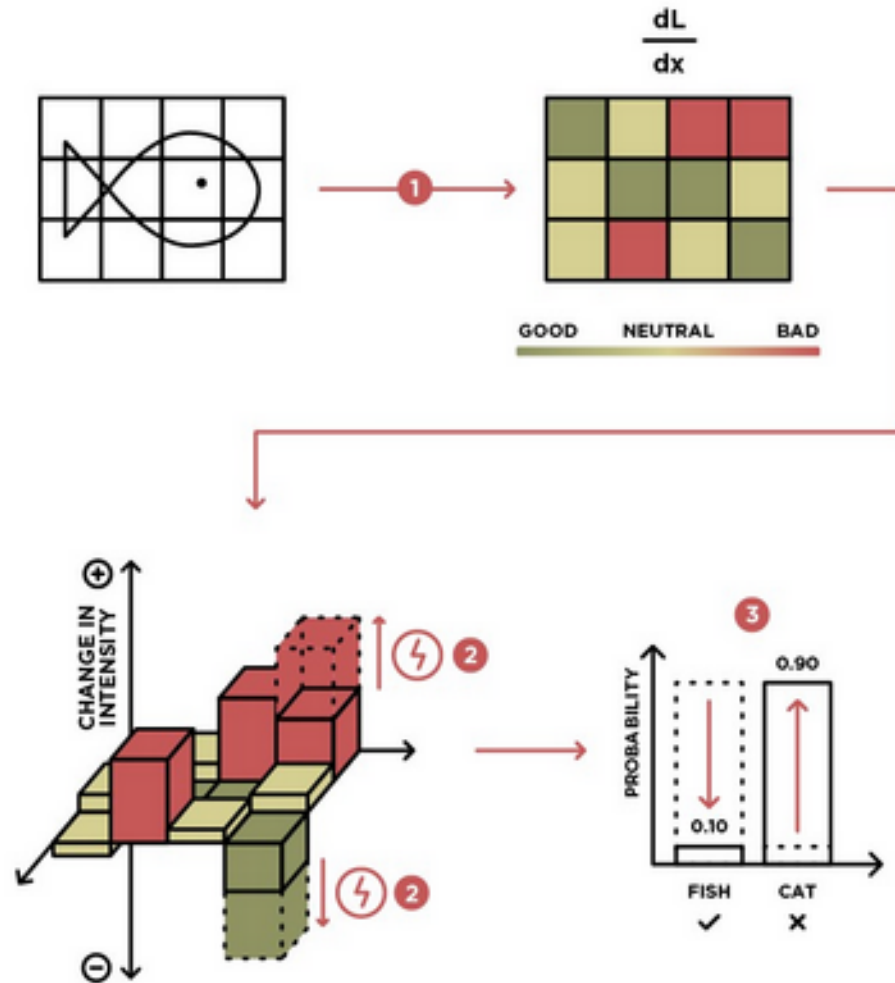
Всё

Black-box

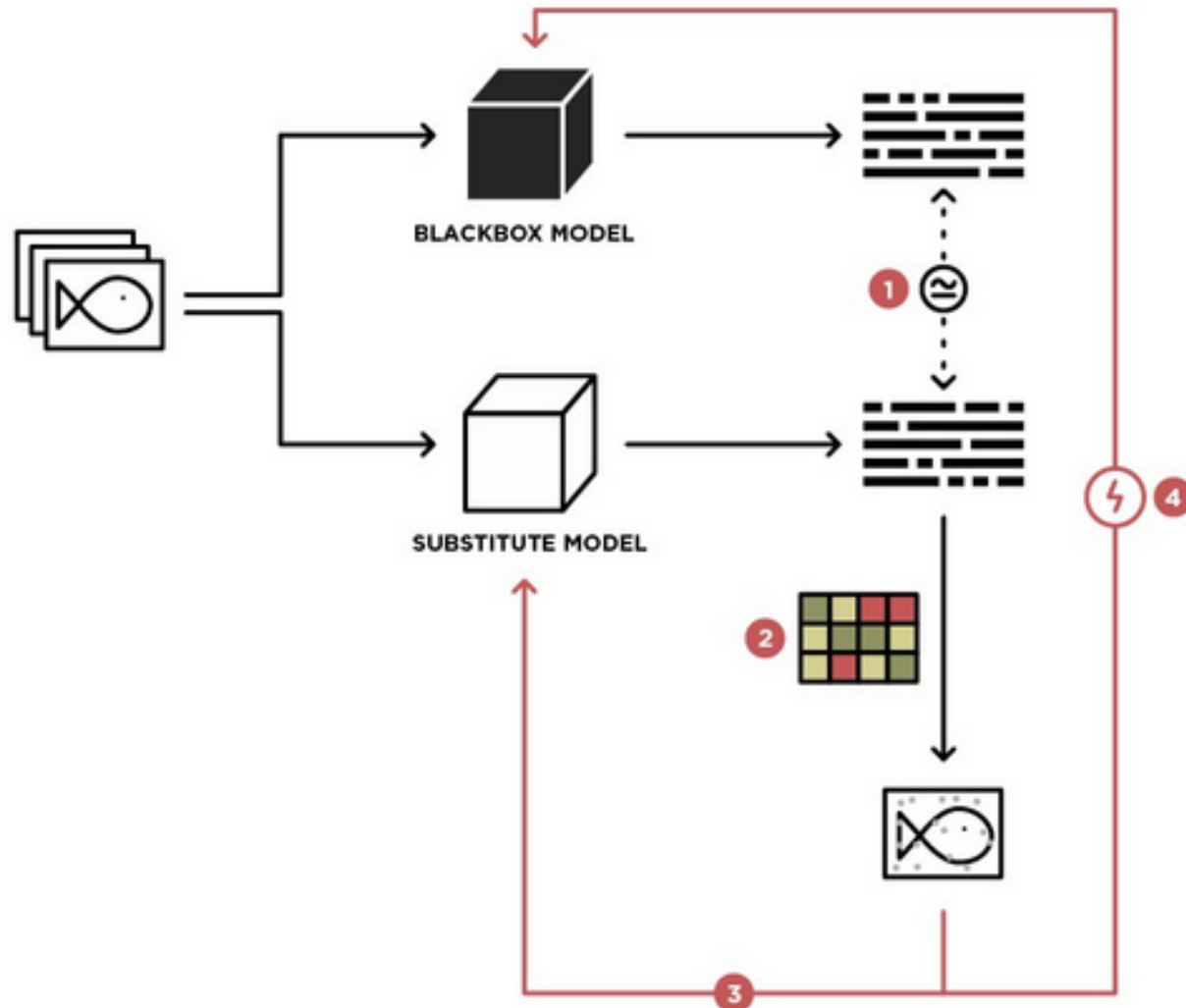


Уверенность

White-box



Black-box



Adversarial Specificity

Targeted attacks

Non-targeted attacks

$$P(f(x') = l') \rightarrow \max$$

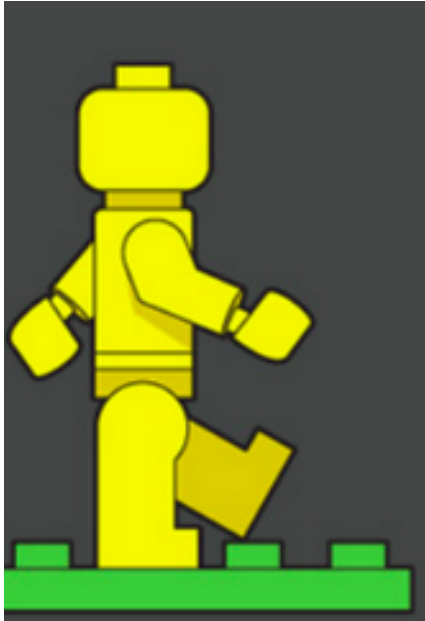
$$\text{a) } P(f(x) = l) \rightarrow \min$$

$$\text{b) } l_k = \min(\eta), l_k \in l_1 \dots l_t$$

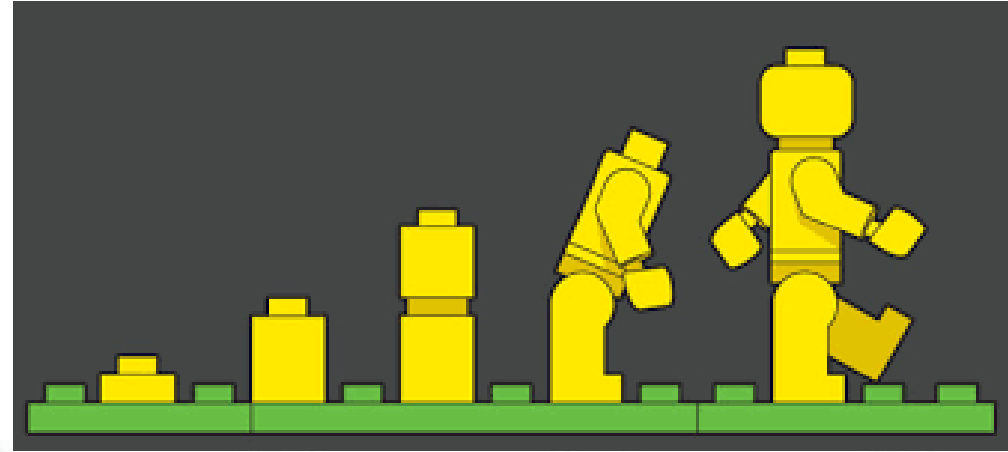
Обозначение	Значение
x	оригинальные данные
l	номер класса ($1 \dots m$)
x'	adversarial example
l'	номер класса (если adv. ex. направ.)
$f(\cdot)$	модель ($f \in F : \mathbb{R}^n \rightarrow l$)
θ	параметры модели f
$J_f(\cdot, \cdot)$	функция потерь для модели f
η	разница между оригиналом и adv. ex. ($\eta = x' - x$)
$\ \cdot\ _p$	l_p норма
∇_t	градиент по t

Attack Frequency

One-time attacks



Iterative attacks



FEEDBACK



II. Как создавать adversarial examples?

L-BFGS

Limited-memory Broyden–Fletcher–Goldfarb–Shanno

Обозначение	Значение
x	оригинальные данные
l	номер класса ($1 \dots m$)
x'	adversarial example
l'	номер класса (если adv. ex. направ.)
$f(\cdot)$	модель ($f \in F : \mathbb{R}^n \rightarrow l$)
θ	параметры модели f
$J_f(\cdot, \cdot)$	функция потерь для модели f
η	разница между оригиналом и adv. ex. ($\eta = x' - x$)
$\ \cdot\ _p$	l_p норма
∇_t	градиент по t

$$\begin{aligned} \min_{x'} \quad & c\|\eta\| + J_{\theta}(x', l') \\ \text{s.t.} \quad & x' \in [0, 1]. \end{aligned}$$

\Rightarrow Лине́йный поиск по c



Страус

FGSM

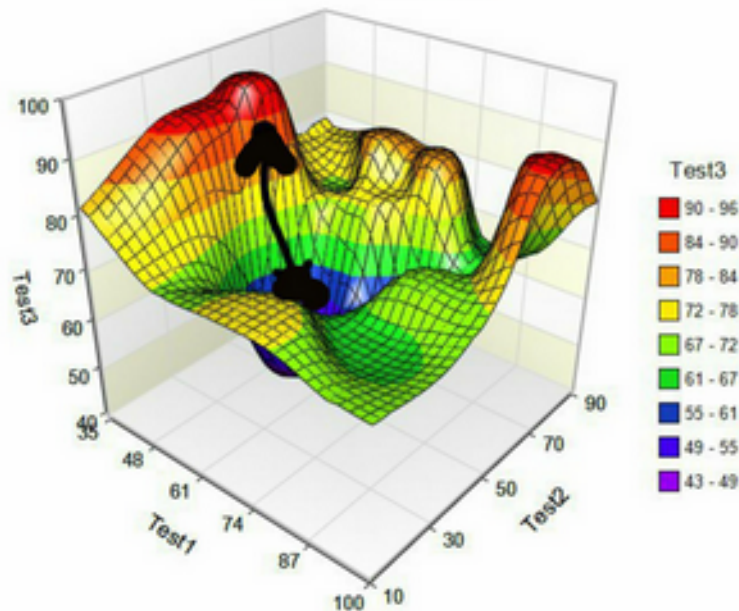
Fast Gradient Sign Method

We are only interested in the sign of the slopes to know if we want to increase or decrease the pixel values

$$\eta = \epsilon \text{sign}(\nabla_x J_\theta(x, l)),$$
$$x' = x + \eta$$

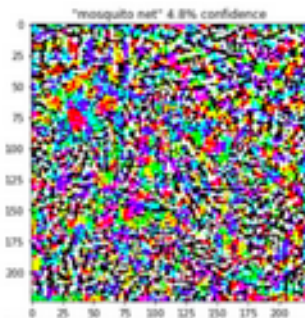
Обозначение	Значение
x	оригинальные данные
l	номер класса ($1 \dots m$)
x'	adversarial example
l'	номер класса (если adv. ex. направ.)
$f(\cdot)$	модель ($f \in F : \mathbb{R}^n \rightarrow l$)
θ	параметры модели f
$J_f(\cdot, \cdot)$	функция потерь для модели f
η	разница между оригиналом и adv. ex. ($\eta = x' - x$)
$\ \cdot\ _p$	l_p норма
∇_t	градиент по t

Surface Plot of Test3



FGSM

Fast Gradient Sign Method

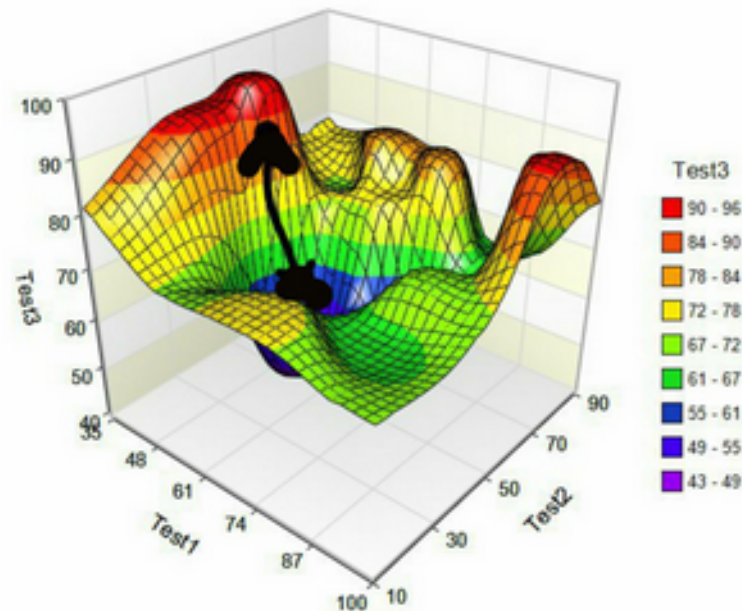


$$\eta = \epsilon \text{sign}(\nabla_x J_\theta(x, l)),$$

$$x' = x + \eta$$

Обозначение	Значение
x	оригинальные данные
l	номер класса ($1 \dots m$)
x'	adversarial example
l'	номер класса (если adv. ex. направ.)
$f(\cdot)$	модель ($f \in F : \mathbb{R}^n \rightarrow l$)
θ	параметры модели f
$J_f(\cdot, \cdot)$	функция потерь для модели f
η	разница между оригиналом и adv. ex. ($\eta = x' - x$)
$\ \cdot\ _p$	l_p норма
∇_t	градиент по t

Surface Plot of Test3



FGV

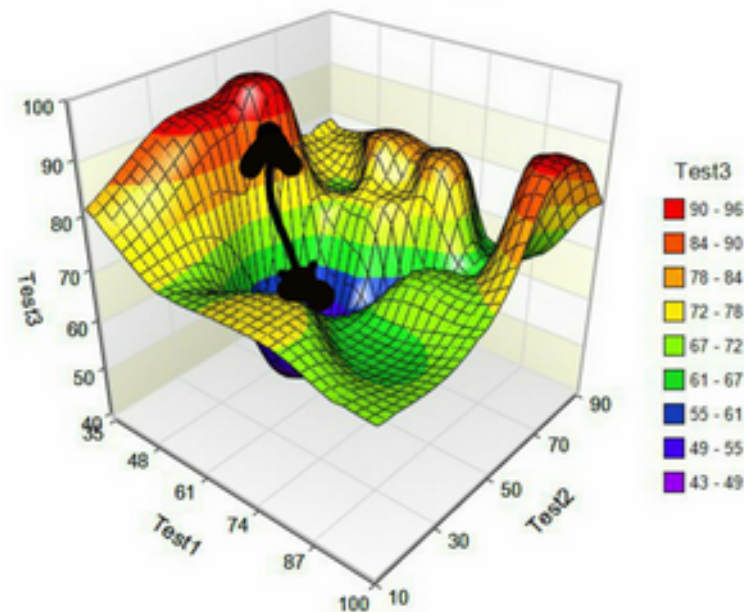
Fast Gradient Value Method

Обозначение	Значение
x	оригинальные данные
l	номер класса ($1 \dots m$)
x'	adversarial example
l'	номер класса (если adv. ex. направ.)
$f(\cdot)$	модель ($f \in F : \mathbb{R}^n \rightarrow l$)
θ	параметры модели f
$J_f(\cdot, \cdot)$	функция потерь для модели f
η	разница между оригиналом и adv. ex. ($\eta = x' - x$)
$\ \cdot\ _p$	l_p норма
∇_t	градиент по t

$$\eta = \nabla_x J(\theta, x, l)$$

$$x' = x + \eta$$

Surface Plot of Test3



FGSM

$$\eta = \epsilon \text{sign}(\nabla_x J_\theta(x, l)),$$
$$x' = x + \eta$$

$$\mathbf{g}_{t+1} = \mu \mathbf{g}_t + \frac{\nabla_x J_\theta(x'_t, l)}{\|\nabla_x J_\theta(x'_t, l)\|},$$

FGSM + Momentum

$$X_0^{adv} = X, \quad X_{N+1}^{adv} = \text{Clip}_{X, \epsilon} \left\{ X_N^{adv} + \alpha \text{sign}(\nabla_X J(X_N^{adv}, y_{true})) \right\}$$

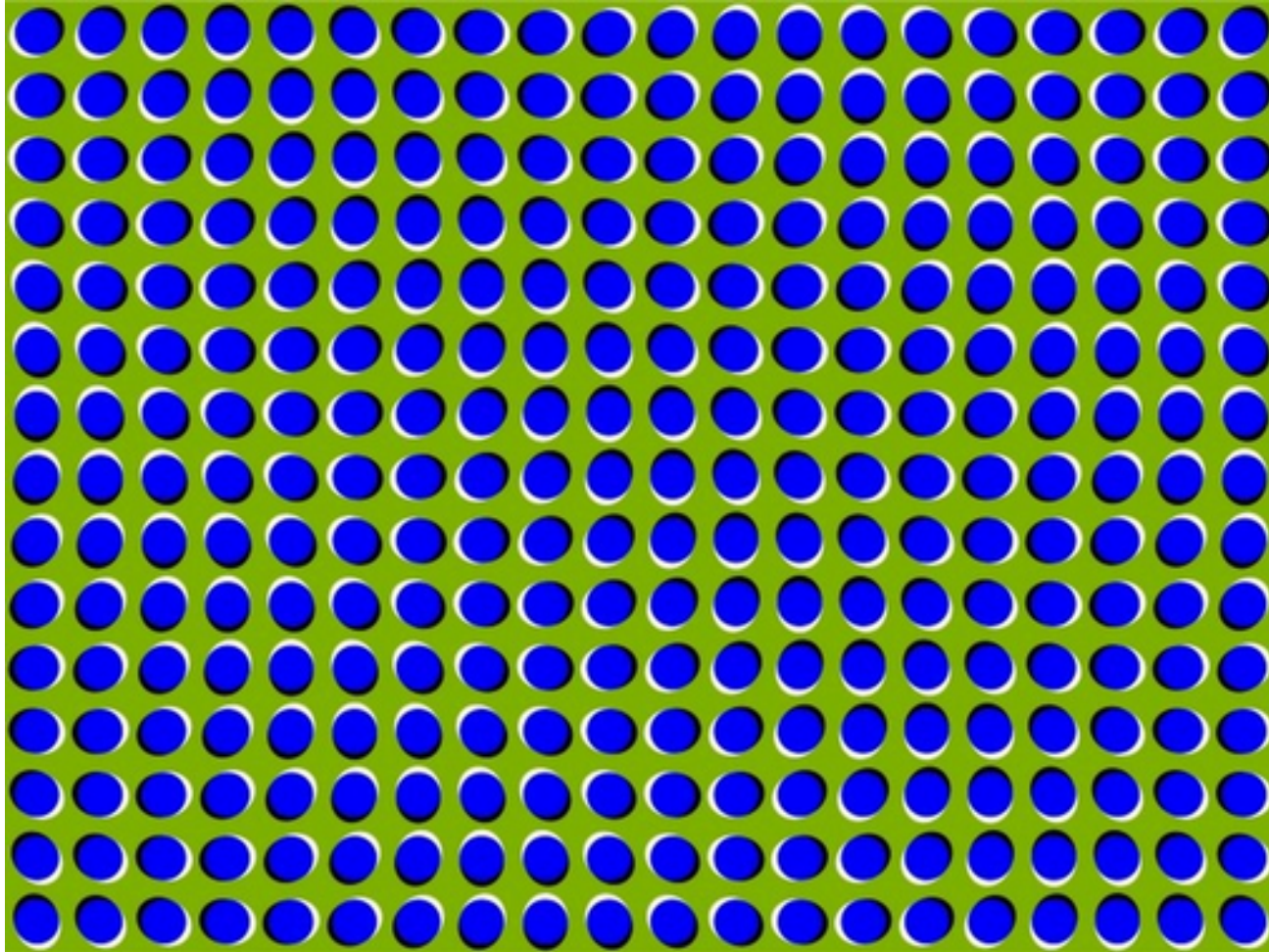
Basic Iterative Method (BIM)

$$x' = x - \epsilon \text{sign}(\nabla_x J(\theta, x, l')).$$

One-step Target Class Method (OTCM)

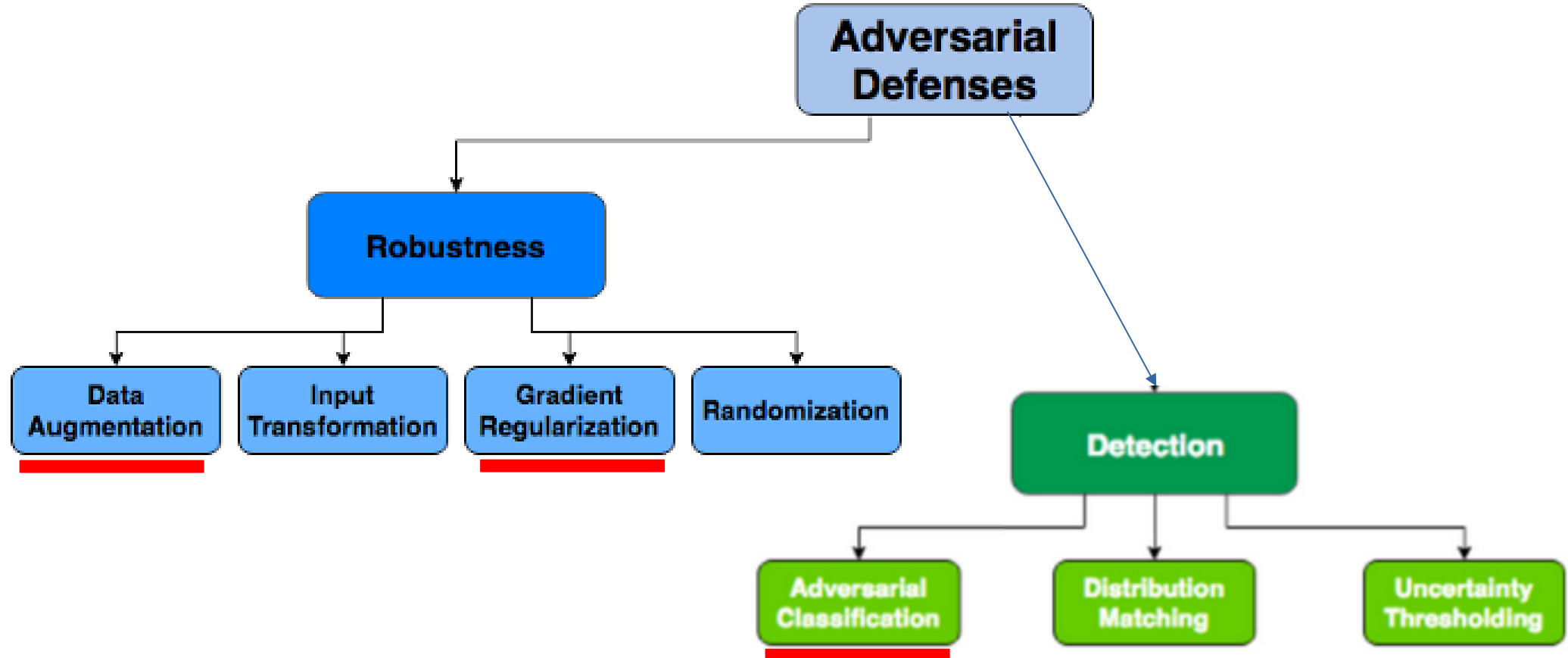
Обозначение	Значение
x	оригинальные данные
l	номер класса ($1 \dots m$)
x'	adversarial example
l'	номер класса (если adv. ex. направ.)
$f(\cdot)$	модель ($f \in F : \mathbb{R}^n \rightarrow l$)
θ	параметры модели f
$J_f(\cdot, \cdot)$	функция потерь для модели f
η	разница между оригиналом и adv. ex. ($\eta = x' - x$)
$\ \cdot\ _p$	l_p норма
∇_t	градиент по t

Почувствуй себя машиной

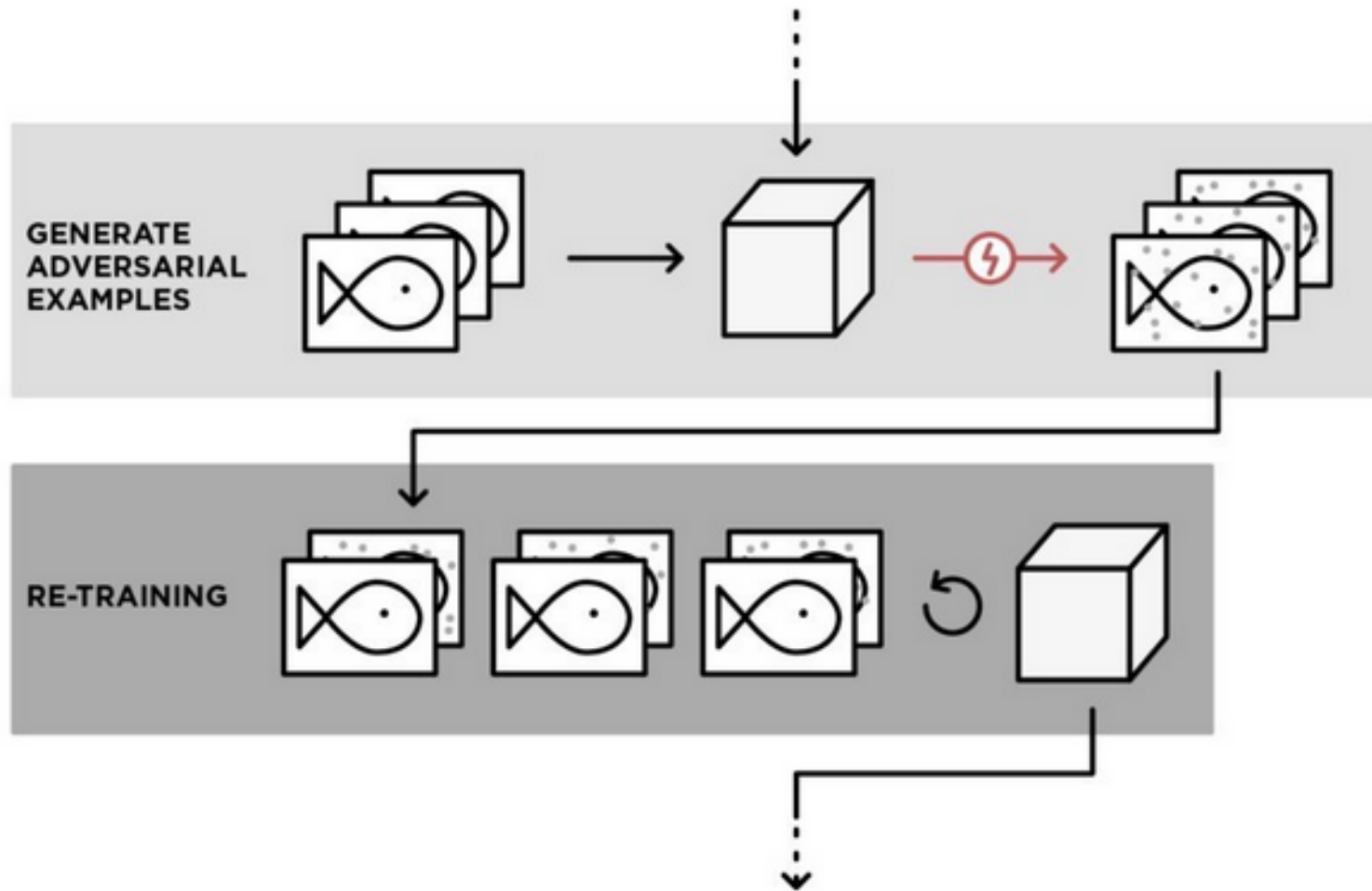


III. Как защищаться от adversarial атак?

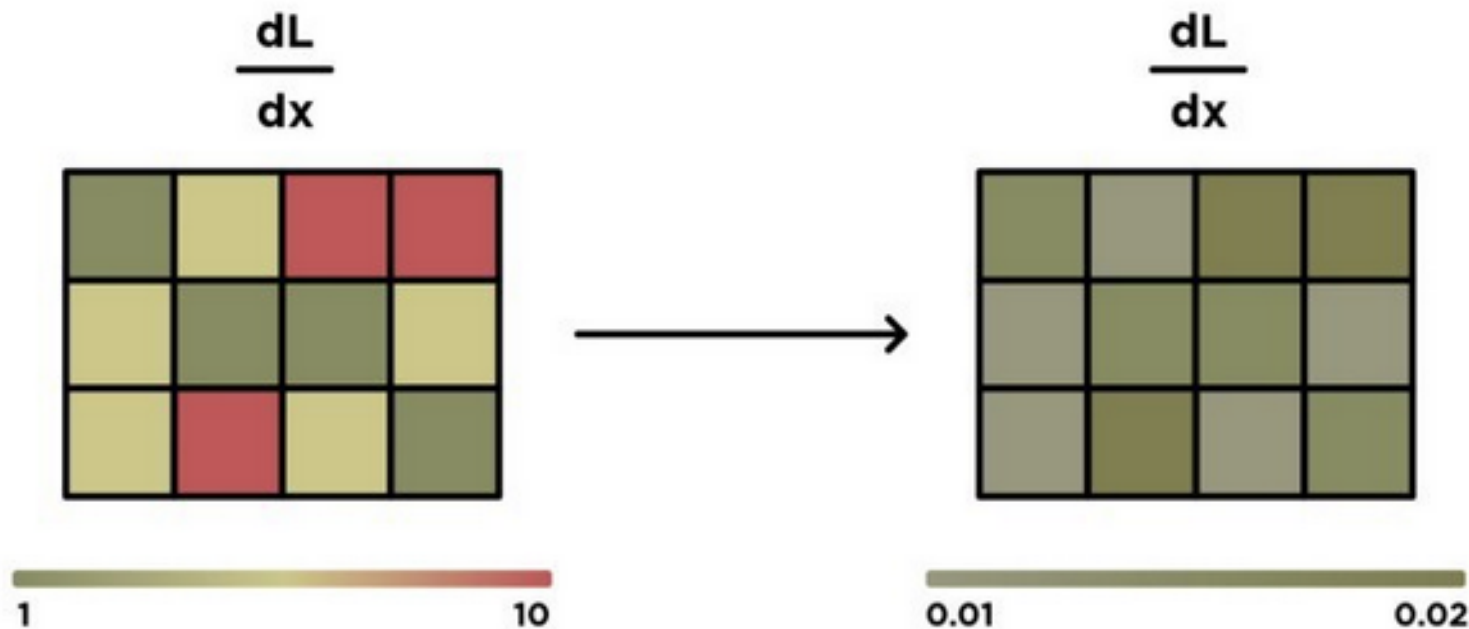
Adversarial Defences



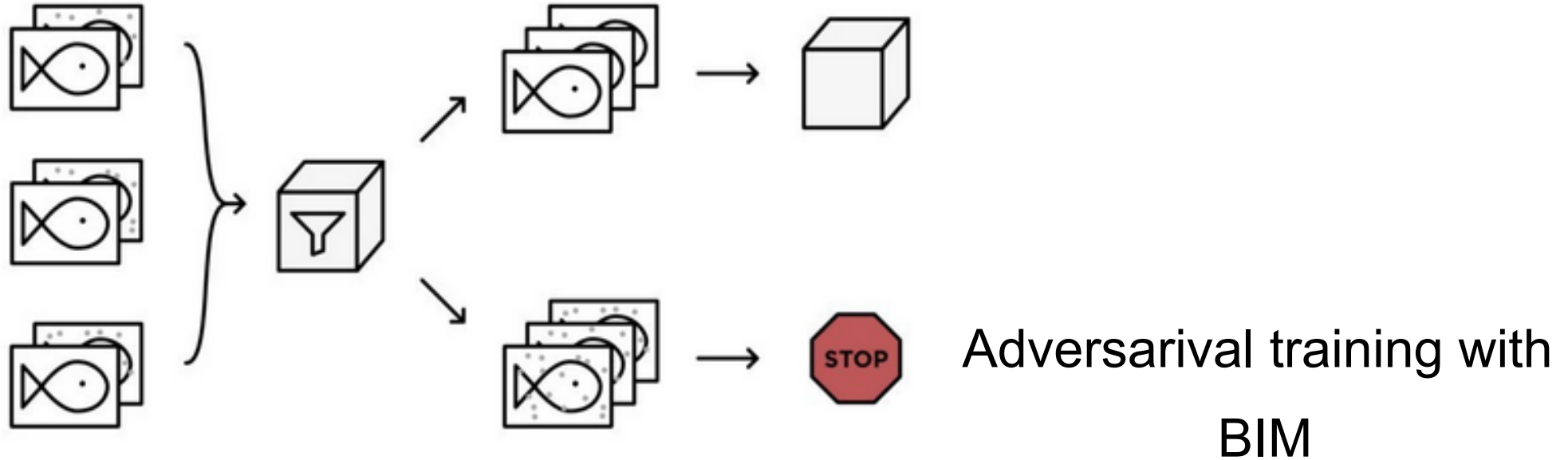
Adversarial Training



Regularizing the Gradient of the Model



Detecting Adversaries Through Classification



$$\mathbf{X}_0^{adv} = \mathbf{X}, \quad \mathbf{X}_{N+1}^{adv} = \text{Clip}_{\mathbf{X}, \epsilon} \left\{ \mathbf{X}_N^{adv} + \alpha \text{sign}(\nabla_{\mathbf{X}} J(\mathbf{X}_N^{adv}, y_{true})) \right\}$$

Face recognition

Input: Hugh Laurie with adversarial glasses



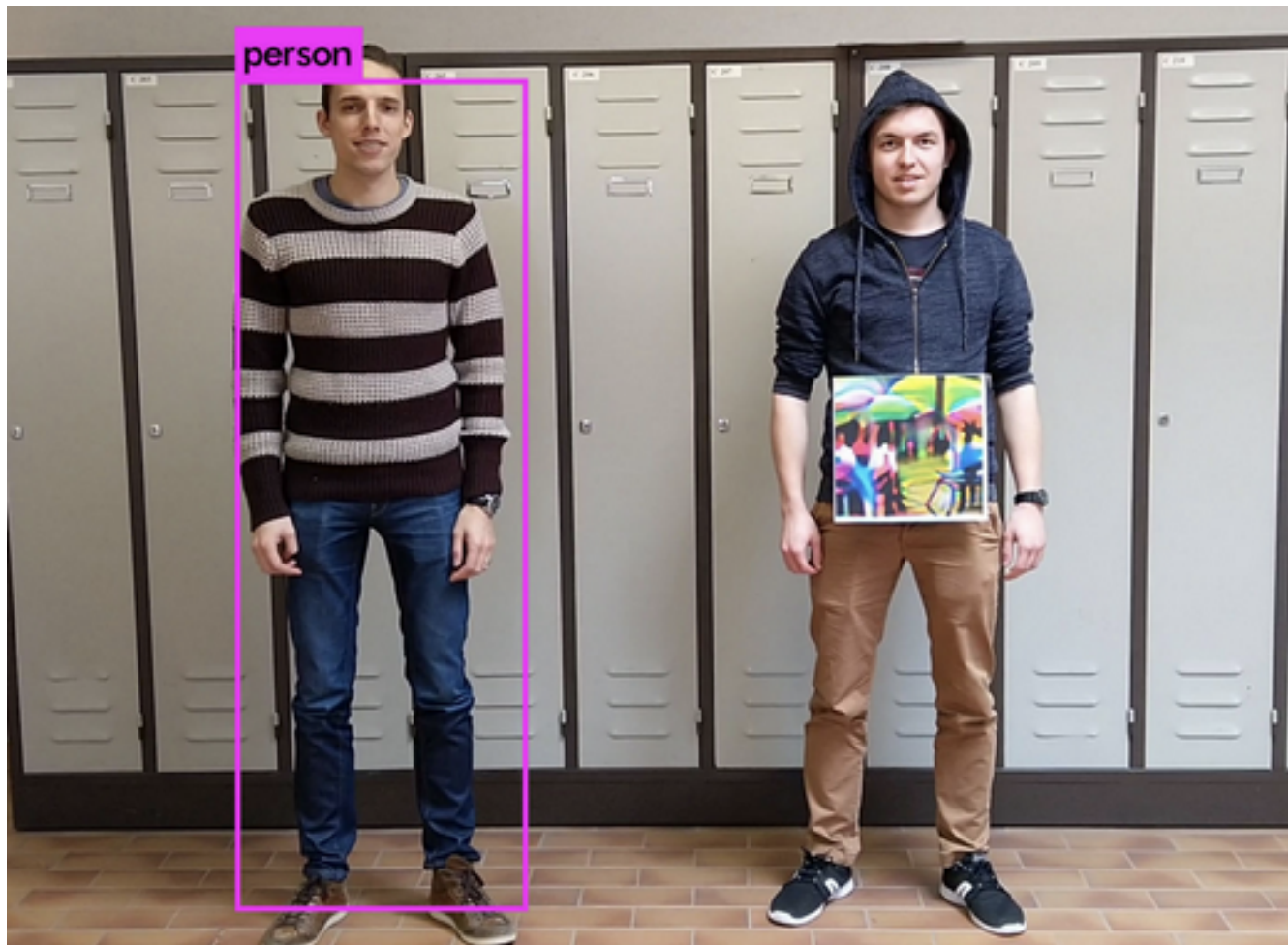
Prediction: Viggo Mortensen



Impersonation Attack



Object detection



WHO WOULD WIN?



**STATE OF THE ART
NEURAL NETWORK**



ONE NOISY BOI

Вопросы:

- Чем отличаются white и black box атаки?
- Как работает метод FGSM?
- Опишите своими словами любой метод Adversarial Defence.

References:

<https://arxiv.org/pdf/1312.6199.pdf>

<https://arxiv.org/pdf/1412.6572.pdf>

<https://arxiv.org/pdf/1712.07107.pdf>

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

Картинки:

<https://medium.com/element-ai-research-lab/securing-machine-learning-models-against-adversarial-attacks-b6cd5d2be8e2>

<https://medium.com/@ml.at.berkeley/tricking-neural-networks-create-your-own-adversarial-examples-a61eb7620fd8>

<https://medium.com/element-ai-research-lab/tricking-a-machine-into-thinking-youre-milla-jovovich-b19bf322d55c>

<https://towardsdatascience.com/adversarial-examples-in-deep-learning-be0b08a94953>

<https://towardsdatascience.com/spam-detection-with-logistic-regression-23e3709e522?gi=324a60138bcc>