

# Применение состязательного обучения к исследованию устойчивости нейросетей в реальном мире

Петюшко А. А.  
petyushko.alexander1@huawei.com

МГУ им. М.В.Ломоносова, механико-математический факультет, кафедра МатИС  
Huawei, Intelligence Systems and Data Science Technnology Center

26 июня 2020 г.

## COMPUTER VISION

SpiceIT Networking



## 1 Intelligence Systems and Data Science Technology Center



- 1 Intelligence Systems and Data Science Technology Center
- 2 Потрясающие успехи СНС в компьютерном зрении



- 1 Intelligence Systems and Data Science Technology Center
- 2 Потрясающие успехи СНС в компьютерном зрении
- 3 (He) устойчивость СНС в компьютерном зрении





- 1 Intelligence Systems and Data Science Technology Center
- 2 Потрясающие успехи СНС в компьютерном зрении
- 3 (Не) устойчивость СНС в компьютерном зрении
- 4 Методы состязательных атак в цифровой области

- 1 Intelligence Systems and Data Science Technology Center
- 2 Потрясающие успехи СНС в компьютерном зрении
- 3 (Не) устойчивость СНС в компьютерном зрении
- 4 Методы состязательных атак в цифровой области
- 5 Методы состязательных атак в реальном мире



- 1 Intelligence Systems and Data Science Technology Center
- 2 Потрясающие успехи СНС в компьютерном зрении
- 3 (Не) устойчивость СНС в компьютерном зрении
- 4 Методы состязательных атак в цифровой области
- 5 Методы состязательных атак в реальном мире
- 6 Состязательные атаки на системы детекции и распознавания лиц в реальном мире



# Intelligence Systems and Data Science Technology Center: научное сотрудничество



ISDSTC



**Skoltech**

Skolkovo Institute of Science and Technology



Санкт-Петербургский  
государственный университет



МОСКОВСКИЙ  
ГОСУДАРСТВЕННЫЙ  
УНИВЕРСИТЕТ ИМЕНИ  
М.В.ЛОМОНОСОВА



В 2019 году в Huawei стартовала образовательная программа **SHARE**: Школа опережающего научного образования Хуавэй (School of Huawei Advanced Research Education).



В 2019 году в Huawei стартовала образовательная программа **SHARE**: Школа опережающего научного образования Хуавэй (School of Huawei Advanced Research Education).

Intelligence Systems and Data Science Technology Center проводит занятия в МГУ им. М.В. Ломоносова:



В 2019 году в Huawei стартовала образовательная программа **SHARE**: Школа опережающего научного образования Хуавэй (School of Huawei Advanced Research Education).

Intelligence Systems and Data Science Technology Center проводит занятия в МГУ им. М.В. Ломоносова:

- 2 года длится программа;



В 2019 году в Huawei стартовала образовательная программа **SHARE**: Школа опережающего научного образования Хуавэй (School of Huawei Advanced Research Education).

Intelligence Systems and Data Science Technology Center проводит занятия в МГУ им. М.В. Ломоносова:

- 2 года длится программа;
- 12 полусеместровых курсов;





В 2019 году в Huawei стартовала образовательная программа **SHARE**: Школа опережающего научного образования Хуавэй (School of Huawei Advanced Research Education).

Intelligence Systems and Data Science Technology Center проводит занятия в МГУ им. М.В. Ломоносова:

- 2 года длится программа;
- 12 полусеместровых курсов;
- 2 направления:



В 2019 году в Huawei стартовала образовательная программа **SHARE**: Школа опережающего научного образования Хуавэй (School of Huawei Advanced Research Education).

Intelligence Systems and Data Science Technology Center проводит занятия в МГУ им. М.В. Ломоносова:

- 2 года длится программа;
- 12 полусеместровых курсов;
- 2 направления:
  - Специализация “Компьютерное зрение и машинное обучение”;



В 2019 году в Huawei стартовала образовательная программа **SHARE**: Школа опережающего научного образования Хуавэй (School of Huawei Advanced Research Education).

Intelligence Systems and Data Science Technology Center проводит занятия в МГУ им. М.В. Ломоносова:

- 2 года длится программа;
- 12 полусеместровых курсов;
- 2 направления:
  - Специализация “Компьютерное зрение и машинное обучение”;
  - Специализация “Большие данные и теория информации”.



Начиная с 1943 года, когда впервые была предложена математическая формализация МакКаломом и Питтсом понятия “искусственного нейрона”, нейросети становились<sup>1</sup>:

---

<sup>1</sup>Image credits: <https://arxiv.org/abs/1409.4842>

Начиная с 1943 года, когда впервые была предложена математическая формализация МакКаломом и Питтсом понятия “искусственного нейрона”, нейросети становились<sup>1</sup>:

- Объемнее (содержали больше параметров),

---

<sup>1</sup>Image credits: <https://arxiv.org/abs/1409.4842>

Начиная с 1943 года, когда впервые была предложена математическая формализация МакКаломом и Питтсом понятия “искусственного нейрона”, нейросети становились<sup>1</sup>:

- Объемнее (содержали больше параметров),
- Глубже (содержали больше блоков вычислений),

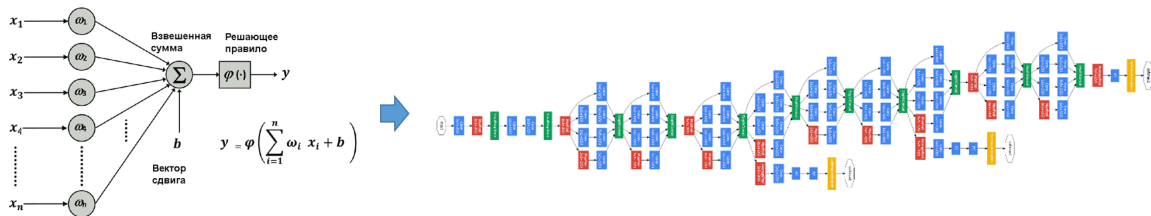
---

<sup>1</sup>Image credits: <https://arxiv.org/abs/1409.4842>

# Развитие нейросетей

Начиная с 1943 года, когда впервые была предложена математическая формализация МакКаломом и Питтсом понятия “искусственного нейрона”, нейросети становились<sup>1</sup>:

- Объемнее (содержали больше параметров),
- Глубже (содержали больше блоков вычислений),
- Лучше! (более правильно решали поставленные перед ними задачи)



<sup>1</sup>Image credits: <https://arxiv.org/abs/1409.4842>

- Для работы с фотографиями и видео лучше всего подходят сверточные нейронные сети (СНС),

---

<sup>2</sup>Image credits: <https://adeshpande3.github.io/>, <https://stepupanalytics.com>



# Сверточные нейросети<sup>2</sup>

- Для работы с фотографиями и видео лучше всего подходят сверточные нейронные сети (СНС),
- Например, позволяют выделять объекты и определять их класс,

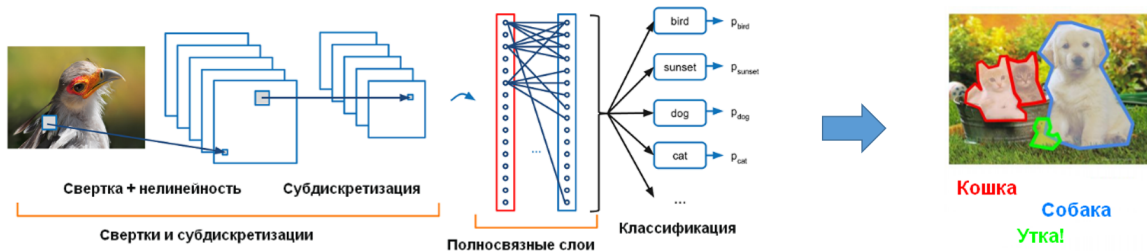
---

<sup>2</sup>Image credits: <https://adeshpande3.github.io/>, <https://stepupanalytics.com>



# Сверточные нейросети<sup>2</sup>

- Для работы с фотографиями и видео лучше всего подходят сверточные нейронные сети (СНС),
- Например, позволяют выделять объекты и определять их класс,
- Ну и отвечают на главный вопрос – кошка или собака?



<sup>2</sup>Image credits: <https://adeshpande3.github.io/>, <https://stepupanalytics.com>

Давайте разберемся, так ли уж хороши сверточные нейросети, действительно ли оправдано все то внимание, которое им уделяют?

---

<sup>3</sup>Image credit: <https://spectrum.ieee.org>

Давайте разберемся, так ли уж хороши сверточные нейросети, действительно ли оправдано все то внимание, которое им уделяют?

## Вопрос1

Как сейчас соотносится качество распознавания человеком и СНС для известных баз данных?

<sup>3</sup>Image credit: <https://spectrum.ieee.org>

Давайте разберемся, так ли уж хороши сверточные нейросети, действительно ли оправдано все то внимание, которое им уделяют?

## Вопрос1

Как сейчас соотносится качество распознавания человеком и СНС для известных баз данных?

## Вопрос2

Насколько устойчивы СНС по отношению к входным данным? Легко ли их сломать?

<sup>3</sup>Image credit: <https://spectrum.ieee.org>

Давайте разберемся, так ли уж хороши сверточные нейросети, действительно ли оправдано все то внимание, которое им уделяют?

## Вопрос1

Как сейчас соотносится качество распознавания человеком и СНС для известных баз данных?

## Вопрос2

Насколько устойчивы СНС по отношению к входным данным? Легко ли их сломать?

CNN vs Human<sup>3</sup>



<sup>3</sup>Image credit: <https://spectrum.ieee.org>

# Человек или СНС?

## ImageNet<sup>4</sup> (1000-классовая база данных изображений)

- Тор-5 ошибка для человека<sup>5</sup>: 5.1%
- Тор-5 ошибка для СНС<sup>6</sup>: 2.0%

<sup>4</sup><http://www.image-net.org/>

<sup>5</sup><http://karpathy.github.io/2014/09/02/>

[what-i-learned-from-competing-against-a-convnet-on-imagenet/](http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/)

<sup>6</sup>Touvron, Hugo, et al. "Fixing the train-test resolution discrepancy." 2019

<sup>7</sup><http://vis-www.cs.umass.edu/lfw/>

<sup>8</sup>Kumar, Neeraj, et al. "Attribute and simile classifiers for face verification." 2009

<sup>9</sup>Deng, Jiankang, et al. "Arcface: Additive angular margin loss for deep face recognition." 2018



# Человек или СНС?

## ImageNet<sup>4</sup> (1000-классовая база данных изображений)

- Тор-5 ошибка для человека<sup>5</sup>: 5.1%
- Тор-5 ошибка для СНС<sup>6</sup>: 2.0%

## Labeled Faces in the Wild<sup>7</sup> (база данных лиц)

- Ошибка верификации для человека<sup>8</sup>: 2.47%
- Ошибка верификации для СНС<sup>9</sup>: 0.17%

<sup>4</sup><http://www.image-net.org/>

<sup>5</sup><http://karpathy.github.io/2014/09/02/>

[what-i-learned-from-competing-against-a-convnet-on-imagenet/](http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/)

<sup>6</sup>Touvron, Hugo, et al. "Fixing the train-test resolution discrepancy." 2019

<sup>7</sup><http://vis-www.cs.umass.edu/lfw/>

<sup>8</sup>Kumar, Neeraj, et al. "Attribute and simile classifiers for face verification." 2009

<sup>9</sup>Deng, Jiankang, et al. "Arcface: Additive angular margin loss for deep face recognition." 2018





# Человек или СНС?

## ImageNet<sup>4</sup> (1000-классовая база данных изображений)

- Тор-5 ошибка для человека<sup>5</sup>: 5.1%
- Тор-5 ошибка для СНС<sup>6</sup>: 2.0%

## Labeled Faces in the Wild<sup>7</sup> (база данных лиц)

- Ошибка верификации для человека<sup>8</sup>: 2.47%
- Ошибка верификации для СНС<sup>9</sup>: 0.17%

<sup>4</sup><http://www.image-net.org/>

<sup>5</sup><http://karpathy.github.io/2014/09/02/>

[what-i-learned-from-competing-against-a-convnet-on-imagenet/](http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/)

<sup>6</sup>Touvron, Hugo, et al. "Fixing the train-test resolution discrepancy." 2019

<sup>7</sup><http://vis-www.cs.umass.edu/lfw/>

<sup>8</sup>Kumar, Neeraj, et al. "Attribute and simile classifiers for face verification." 2009

<sup>9</sup>Deng, Jiankang, et al. "Arcface: Additive angular margin loss for deep face recognition." 2018

LFW



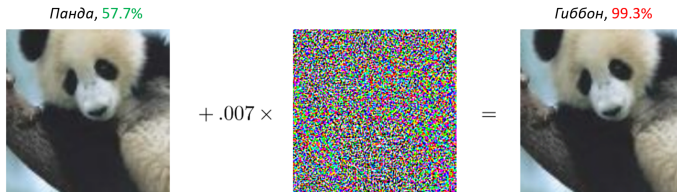
- Можно внести практически незаметные для глаза человека возмущения во входные данные, которые, тем не менее, полностью поменяют выход нейронной сети

---

<sup>10</sup>Image credit: <https://arxiv.org/pdf/1412.6572.pdf>

# Такие неустойчивые СНС

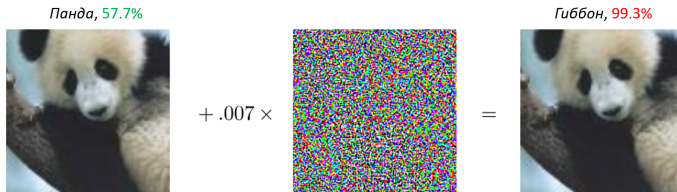
- Можно внести практически незаметные для глаза человека возмущения во входные данные, которые, тем не менее, полностью поменяют выход нейронной сети
- Например, результат классификации с “панды” поменяется на “гиббона”<sup>10</sup>



<sup>10</sup>Image credit: <https://arxiv.org/pdf/1412.6572.pdf>

# Такие неустойчивые СНС

- Можно внести практически незаметные для глаза человека возмущения во входные данные, которые, тем не менее, полностью поменяют выход нейронной сети
- Например, результат классификации с “панды” поменяется на “гиббона”<sup>10</sup>



Такое возмущение называется **сопоставительной атакой** (adversarial attack)

<sup>10</sup>Image credit: <https://arxiv.org/pdf/1412.6572.pdf>

# Атака СНС, предназначенных для сегментации или обнаружения

- Можно атаковать также СНС, которые не предназначены для классификации — например, для обнаружения и сегментации изображений<sup>11</sup>

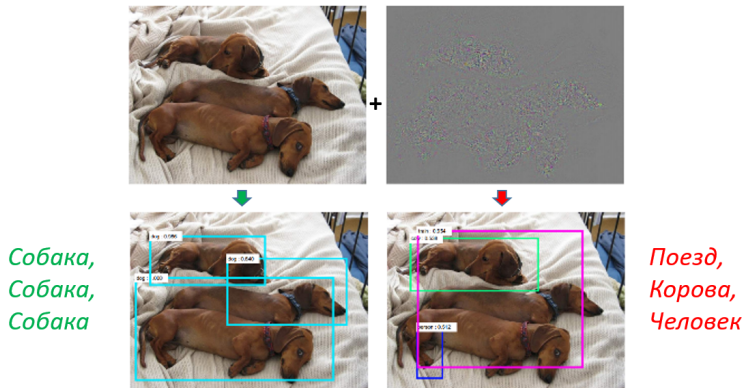
---

<sup>11</sup>Xie, Cihang, et al. "Adversarial examples for semantic segmentation and object detection." 2017. ▶



# Атака СНС, предназначенных для сегментации или обнаружения

- Можно атаковать также СНС, которые не предназначены для классификации — например, для обнаружения и сегментации изображений<sup>11</sup>



<sup>11</sup>Xie, Cihang, et al. "Adversarial examples for semantic segmentation and object detection." 2017.

# Атака нейросетей, не предназначенных для изображений

- Можно атаковать даже НС, которые вообще не работают с изображениями — например, НС для вопросно-ответных систем (QA, question answering systems)<sup>12</sup>

---

<sup>12</sup> Jia, Robin, and Percy Liang. "Adversarial examples for evaluating reading comprehension systems." 2017



# Атака нейросетей, не предназначенных для изображений

- Можно атаковать даже НС, которые вообще не работают с изображениями — например, НС для вопросно-ответных систем (QA, question answering systems)<sup>12</sup>

**Article:** Super Bowl 50

**Paragraph:** *“Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver’s Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV.”*

**Question:** *“What is the name of the quarterback who was 38 in Super Bowl XXXIII?”*

**Original Prediction:** John Elway

**Prediction under adversary:** Jeff Dean

<sup>12</sup>Jia, Robin, and Percy Liang. “Adversarial examples for evaluating reading comprehension systems.” 2017



# Одна из главных причин существования атак

- Одна из основных причин такого поведения СНС на похожих изображениях — неустойчивость СНС

---

<sup>13</sup>Image credit: <https://secml.github.io/>

<sup>14</sup>Fawzi, Alhussein, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. "Robustness of classifiers: from adversarial to random noise." 2016

# Одна из главных причин существования атак

- Одна из основных причин такого поведения СНС на похожих изображениях — неустойчивость СНС
- А именно, разделяющие границы классификатора часто проходят очень близко к обучающим данным, и легко “заступить” за такую границу<sup>13,14</sup>

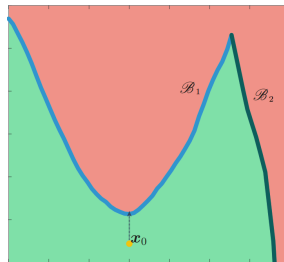
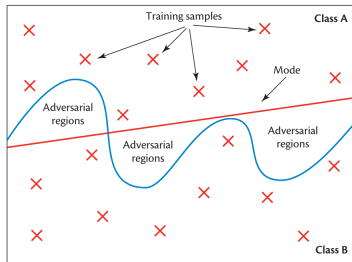
---

<sup>13</sup>Image credit: <https://secml.github.io/>

<sup>14</sup>Fawzi, Alhussein, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. “Robustness of classifiers: from adversarial to random noise.” 2016

# Одна из главных причин существования атак

- Одна из основных причин такого поведения СНС на похожих изображениях — неустойчивость СНС
- А именно, разделяющие границы классификатора часто проходят очень близко к обучающим данным, и легко “заступить” за такую границу<sup>13,14</sup>



<sup>13</sup>Image credit: <https://secml.github.io/>

<sup>14</sup>Fawzi, Alhussein, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard. “Robustness of classifiers: from adversarial to random noise.” 2016

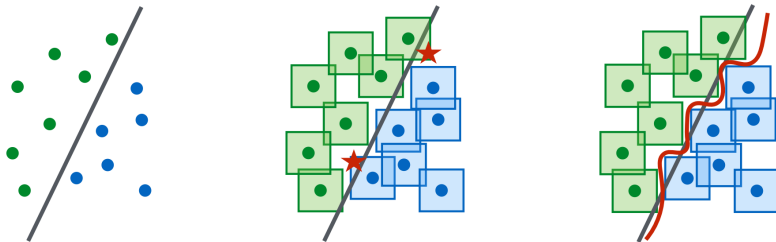
- Поскольку можно обмануть СНС путем небольшого пиксельного возмущения, то почему бы во время обучения для каждого обучающего примера не добавлять и всю его попиксельную окрестность (по некоторой норме, например,  $\ell_\infty$ )<sup>15</sup>

---

<sup>15</sup>Madry, Aleksander, et al. "Towards deep learning models resistant to adversarial attacks." 2017



- Поскольку можно обмануть СНС путем небольшого пиксельного возмущения, то почему бы во время обучения для каждого обучающего примера не добавлять и всю его попиксельную окрестность (по некоторой норме, например,  $\ell_\infty$ )<sup>15</sup>



<sup>15</sup>Madry, Aleksander, et al. "Towards deep learning models resistant to adversarial attacks." 2017

# Простой, но не работающий метод защиты

- Предположим, что исходная картинка размера  $100 \times 100$  пикселей, 3 цвета RGB



# Простой, но не работающий метод защиты

- Предположим, что исходная картинка размера  $100 \times 100$  пикселей, 3 цвета RGB
- Предположим, что наш глаз не сильно различает колебания цвета пикселей в 2 градации (из 256): в каждой точке для каждого цвета можем позволить  $\pm 1$  значение



# Простой, но не работающий метод защиты

- Предположим, что исходная картинка размера  $100 \times 100$  пикселей, 3 цвета RGB
- Предположим, что наш глаз не сильно различает колебания цвета пикселей в 2 градации (из 256): в каждой точке для каждого цвета можем позволить  $\pm 1$  значение
- Тогда для каждого обучающего примера нужно добавить следующее количество его пиксельных соседей:

$$2^{3 \times 100 \times 100} = 2^{30000} = (2^{10})^{3000} \approx (10^3)^{3000} = 10^{9000}$$





# Простой, но не работающий метод защиты

- Предположим, что исходная картинка размера  $100 \times 100$  пикселей, 3 цвета RGB
- Предположим, что наш глаз не сильно различает колебания цвета пикселей в 2 градации (из 256): в каждой точке для каждого цвета можем позволить  $\pm 1$  значение
- Тогда для каждого обучающего примера нужно добавить следующее количество его пиксельных соседей:

$$2^{3 \times 100 \times 100} = 2^{30000} = (2^{10})^{3000} \approx (10^3)^{3000} = 10^{9000}$$

- Это гораздо больше числа атомов в видимой части Вселенной ( $10^{80}$ )!



# Простой, но не работающий метод защиты

- Предположим, что исходная картинка размера  $100 \times 100$  пикселей, 3 цвета RGB
- Предположим, что наш глаз не сильно различает колебания цвета пикселей в 2 градации (из 256): в каждой точке для каждого цвета можем позволить  $\pm 1$  значение
- Тогда для каждого обучающего примера нужно добавить следующее количество его пиксельных соседей:

$$2^{3 \times 100 \times 100} = 2^{30000} = (2^{10})^{3000} \approx (10^3)^{3000} = 10^{9000}$$

- Это гораздо больше числа атомов в видимой части Вселенной ( $10^{80}$ )!
- В общем, не очень реалистично



- Давайте не перебирать всю окрестность обучающего примера, а брать только те точки из окрестности, которые ближе всего к разделяющей поверхности

---

<sup>16</sup>Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." 2014



# Работающий метод защиты

- Давайте не перебирать всю окрестность обучающего примера, а брать только те точки из окрестности, которые ближе всего к разделяющей поверхности
- Такой метод обучения называется состязательным (adversarial training)<sup>16</sup>

---

<sup>16</sup>Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." 2014



# Работающий метод защиты

- Давайте не перебирать всю окрестность обучающего примера, а брать только те точки из окрестности, которые ближе всего к разделяющей поверхности
- Такой метод обучения называется состязательным (adversarial training)<sup>16</sup>

## Плюсы состязательного обучения

- Не нужно перебирать всю окрестность огромной мощности
- В целом, защищает от метода нахождения состязательных примеров

<sup>16</sup>Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." 2014



# Работающий метод защиты

- Давайте не перебирать всю окрестность обучающего примера, а брать только те точки из окрестности, которые ближе всего к разделяющей поверхности
- Такой метод обучения называется состязательным (adversarial training)<sup>16</sup>

## Плюсы состязательного обучения

- Не нужно перебирать всю окрестность огромной мощности
- В целом, защищает от метода нахождения состязательных примеров

## Минусы состязательного обучения

- Процедура нахождения хороших состязательных примеров работает медленно (гораздо медленнее одного градиентного шага)
- Защищает **только** от того метода нахождения состязательных примеров, который использовался в состязательном обучении

<sup>16</sup>Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." 2014



# Состязательные атаки: необходимые обозначения

- Пусть  $x \in B = [0, 1]^{C \times M \times N}$  — входная картинка  $C \times M \times N$ , где  $C$  — количество цветов (1 для ч/б, 3 для RGB)
- $y_{gt}$  — правильный класс для  $x$
- $\theta$  — параметры СНС-классификатора
- $L(\theta, x, y_{gt})$  — функция потерь
- $f(x)$  — выход классификатора (распознанный класс); при обучении мы добиваемся равенства  $f(x) = y_{gt}$



# Состязательные атаки: необходимые обозначения

- Пусть  $x \in B = [0, 1]^{C \times M \times N}$  — входная картинка  $C \times M \times N$ , где  $C$  — количество цветов (1 для ч/б, 3 для RGB)
- $y_{gt}$  — правильный класс для  $x$
- $\theta$  — параметры СНС-классификатора
- $L(\theta, x, y_{gt})$  — функция потерь
- $f(x)$  — выход классификатора (распознанный класс); при обучении мы добиваемся равенства  $f(x) = y_{gt}$
- $r \in B = [0, 1]^{C \times M \times N}$  — аддитивная добавка ко входу  $x$





# Состязательная атака, устойчивость: формулировка

## Цель состязательной атаки

Поменять выход классификатора  $f$  на неправильный путем добавления минимального по некоторой норме (на практике используются  $\ell_0$ ,  $\ell_1$ ,  $\ell_2$  и  $\ell_\infty$  — обозначим через  $\ell_p$ ) возмущения  $r$ , а именно:



## Цель состязательной атаки

Поменять выход классификатора  $f$  на неправильный путем добавления минимального по некоторой норме (на практике используются  $\ell_0$ ,  $\ell_1$ ,  $\ell_2$  и  $\ell_\infty$  — обозначим через  $\ell_p$ ) возмущения  $r$ , а именно: минимизировать  $\|r\|_p$  т.ч.

- 1  $f(x) = y_{gt}$
- 2  $f(x + r) \neq y_{gt}$
- 3  $x + r \in B$



# Состязательная атака, устойчивость: формулировка

## Цель состязательной атаки

Поменять выход классификатора  $f$  на неправильный путем добавления минимального по некоторой норме (на практике используются  $\ell_0$ ,  $\ell_1$ ,  $\ell_2$  и  $\ell_\infty$  — обозначим через  $\ell_p$ ) возмущения  $r$ , а именно: минимизировать  $\|r\|_p$  т.ч.

- 1  $f(x) = y_{gt}$
- 2  $f(x + r) \neq y_{gt}$
- 3  $x + r \in B$

## Устойчивость классификатора

Найти такой класс возмущения  $S(x, f) \subseteq B$ , при котором классификатор не меняет свой выход:

$$f(x + r) = f(x) \quad \forall r \in S(x, f)$$

В обозначениях выше обычное обучение можно сформулировать как

Обучение на примерах

$$\min_{\theta} \mathbb{E}_{x, y_{gt}} [L(\theta, x, y_{gt})]$$



# Состязательное обучение: формулировка

В обозначениях выше обычное обучение можно сформулировать как

## Обучение на примерах

$$\min_{\theta} \mathbb{E}_{x, y_{gt}} [L(\theta, x, y_{gt})]$$

В состязательном обучении мы сначала генерируем (например, каким-нибудь методом атаки) самый сложный пример из некоторой окрестности  $\Delta$  входного примера (например, по  $\ell_p$ -норме), а уже затем минимизируем по параметрам нейросети:



# Состязательное обучение: формулировка

В обозначениях выше обычное обучение можно сформулировать как

## Обучение на примерах

$$\min_{\theta} \mathbb{E}_{x, y_{gt}} [L(\theta, x, y_{gt})]$$

В состязательном обучении мы сначала генерируем (например, каким-нибудь методом атаки) самый сложный пример из некоторой окрестности  $\Delta$  входного примера (например, по  $\ell_p$ -норме), а уже затем минимизируем по параметрам нейросети:

## Состязательное обучение

$$\min_{\theta} \mathbb{E}_{x, y_{gt}} [\max_{r \in \Delta} L(\theta, x + r, y_{gt})]$$



Напомним наиболее употребительные нормы  $\ell_p$  для  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ :



Напомним наиболее употребительные нормы  $\ell_p$  для  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ :

- $\ell_2$ :  $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$





Напомним наиболее употребительные нормы  $\ell_p$  для  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ :

- $\ell_2$ :  $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$
- $\ell_1$ :  $\|x\|_1 = \sum_{i=1}^n |x_i|$



Напомним наиболее употребительные нормы  $\ell_p$  для  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ :

- $\ell_2$ :  $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$
- $\ell_1$ :  $\|x\|_1 = \sum_{i=1}^n |x_i|$
- $\ell_\infty$ :  $\|x\|_\infty = \max_i |x_i|$



Напомним наиболее употребительные нормы  $\ell_p$  для  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ :

- $\ell_2$ :  $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$
- $\ell_1$ :  $\|x\|_1 = \sum_{i=1}^n |x_i|$
- $\ell_\infty$ :  $\|x\|_\infty = \max_i |x_i|$
- $\ell_0$ :  $\|x\|_0 = \sum_{i=1}^n \mathbf{1}_{x_i \neq 0}$



Напомним наиболее употребительные нормы  $\ell_p$  для  $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ :

- $\ell_2$ :  $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$
- $\ell_1$ :  $\|x\|_1 = \sum_{i=1}^n |x_i|$
- $\ell_\infty$ :  $\|x\|_\infty = \max_i |x_i|$
- $\ell_0$ :  $\|x\|_0 = \sum_{i=1}^n \mathbf{1}_{x_i \neq 0}$

**Замечание.** Для  $0 < p < 1$  норма  $\ell_p$ , для которой  $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ , не является нормой



- Изначально устойчивость СНС изучалась с точки зрения адекватной реакции на разные входы

---

<sup>17</sup>Nguyen, Anh, Jason Yosinski, and Jeff Clune. "Deep neural networks are easily fooled: High confidence predictions for unrecognizable images." 2014

# Предтеча состязательных атак

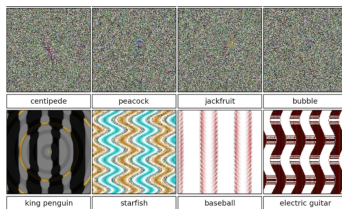
- Изначально устойчивость СНС изучалась с точки зрения адекватной реакции на разные входы
- Выяснилось, что существуют примеры (структурированные или нет), которые на выходе СНС могут давать с большой вероятностью любой класс

---

<sup>17</sup>Nguyen, Anh, Jason Yosinski, and Jeff Clune. "Deep neural networks are easily fooled: High confidence predictions for unrecognizable images." 2014

# Предтеча состязательных атак

- Изначально устойчивость СНС изучалась с точки зрения адекватной реакции на разные входы
- Выяснилось, что существуют примеры (структурированные или нет), которые на выходе СНС могут давать с большой вероятностью любой класс
- Такие примеры назывались “обманными изображениями”<sup>17</sup> (fooling images) и строились с помощью эволюционных алгоритмов



<sup>17</sup>Nguyen, Anh, Jason Yosinski, and Jeff Clune. “Deep neural networks are easily fooled: High confidence predictions for unrecognizable images.” 2014

- **Предложение:** использовать линейную часть функции потерь в окрестности  $x$  и идти по градиенту — FGSM<sup>18</sup> (**F**ast **G**radient **S**ign **M**ethod):

$$r = \epsilon \cdot \text{sign}(\nabla_x L(\theta, x, y_t))$$

где  $0 < \epsilon < 1$  — некоторая константа

---

<sup>18</sup>Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." 2014





- **Предложение:** использовать линейную часть функции потерь в окрестности  $x$  и идти по градиенту — FGSM<sup>18</sup> (**F**ast **G**radient **S**ign **M**ethod):

$$r = \epsilon \cdot \text{sign}(\nabla_x L(\theta, x, y_t))$$

где  $0 < \epsilon < 1$  — некоторая константа

- **Напоминание:** для оптимизации весов СНС применяется метод обратного распространения ошибок, где берется градиент по весам СНС, т.е.  $\nabla_\theta L(\theta, x, y_{gt})$

---

<sup>18</sup>Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." 2014



- **Предложение:** использовать линейную часть функции потерь в окрестности  $x$  и идти по градиенту — FGSM<sup>18</sup> (**F**ast **G**radient **S**ign **M**ethod):

$$r = \epsilon \cdot \text{sign}(\nabla_x L(\theta, x, y_t))$$

где  $0 < \epsilon < 1$  — некоторая константа

- **Напоминание:** для оптимизации весов СНС применяется метод обратного распространения ошибок, где берется градиент по весам СНС, т.е.  $\nabla_{\theta} L(\theta, x, y_{gt})$
- Исследуется норма возмущения  $\ell_{\infty}$  как наиболее близкая к тому, что использует человек

---

<sup>18</sup>Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." 2014



## Метод атаки: I-FGSM (PGD)

- Часто линейная оценка окрестности функции достаточно грубая, и один шаг FGSM порой не приводит к хорошей атаке

---

<sup>19</sup>Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. "Adversarial examples in the physical world." 2016



## Метод атаки: I-FGSM (PGD)

- Часто линейная оценка окрестности функции достаточно грубая, и один шаг FGSM порой не приводит к хорошей атаке
- Для этого применяют итеративный метод I-FGSM<sup>19</sup> (Iterative FGSM), который позволяет двигаться в сторону границы классификатора более точно

---

<sup>19</sup>Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. "Adversarial examples in the physical world." 2016



## Метод атаки: I-FGSM (PGD)

- Часто линейная оценка окрестности функции достаточно грубая, и один шаг FGSM порой не приводит к хорошей атаке
- Для этого применяют итеративный метод I-FGSM<sup>19</sup> (Iterative FGSM), который позволяет двигаться в сторону границы классификатора более точно
- Если  $\Pi_B$  — проекция на  $B$ , то в случае ненаправленной атаки

$$x^{n+1} = \Pi_B(x^n + \text{sign } \nabla_x L(\theta, x, y_{gt})), \quad x^0 = x$$

<sup>19</sup>Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. "Adversarial examples in the physical world." 2016



## Метод атаки: I-FGSM (PGD)

- Часто линейная оценка окрестности функции достаточно грубая, и один шаг FGSM порой не приводит к хорошей атаке
- Для этого применяют итеративный метод I-FGSM<sup>19</sup> (Iterative FGSM), который позволяет двигаться в сторону границы классификатора более точно
- Если  $\Pi_B$  — проекция на  $B$ , то в случае ненаправленной атаки

$$x^{n+1} = \Pi_B(x^n + \text{sign } \nabla_x L(\theta, x, y_{gt})), \quad x^0 = x$$

- Если принять  $\|x - x_{adv}\|_\infty \leq \epsilon$ , то авторы предлагают делать  $n = \min(256\epsilon + 4, 320\epsilon)$  шагов

<sup>19</sup>Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. "Adversarial examples in the physical world." 2016



## Метод атаки: I-FGSM (PGD)

- Часто линейная оценка окрестности функции достаточно грубая, и один шаг FGSM порой не приводит к хорошей атаке
- Для этого применяют итеративный метод I-FGSM<sup>19</sup> (Iterative FGSM), который позволяет двигаться в сторону границы классификатора более точно
- Если  $\Pi_B$  — проекция на  $B$ , то в случае ненаправленной атаки

$$x^{n+1} = \Pi_B(x^n + \text{sign } \nabla_x L(\theta, x, y_{gt})), \quad x^0 = x$$

- Если принять  $\|x - x_{adv}\|_\infty \leq \epsilon$ , то авторы предлагают делать  $n = \min(256\epsilon + 4, 320\epsilon)$  шагов
- Этот метод также называется PGD (**P**rojected **G**radient **D**escent)

<sup>19</sup>Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. "Adversarial examples in the physical world." 2016



- **Замечание:** Методы атаки все больше похожи на шаги оптимизатора

---

<sup>20</sup>Dong, Yinpeng, et al. "Boosting adversarial attacks with momentum." 2017





- **Замечание:** Методы атаки все больше похожи на шаги оптимизатора
- **Идея:** давайте использовать сглаживание градиента — MI-FGSM<sup>20</sup> (Momentum I-FGSM)

---

<sup>20</sup>Dong, Yinpeng, et al. "Boosting adversarial attacks with momentum." 2017



- **Замечание:** Методы атаки все больше похожи на шаги оптимизатора
- **Идея:** давайте использовать сглаживание градиента — MI-FGSM<sup>20</sup> (Momentum I-FGSM)

---

**Algorithm 1** MI-FGSM

---

**Input:** A classifier  $f$  with loss function  $J$ ; a real example  $\mathbf{x}$  and ground-truth label  $y$ ;

**Input:** The size of perturbation  $\epsilon$ ; iterations  $T$  and decay factor  $\mu$ .

**Output:** An adversarial example  $\mathbf{x}^*$  with  $\|\mathbf{x}^* - \mathbf{x}\|_\infty \leq \epsilon$ .

- 1:  $\alpha = \epsilon/T$ ;
- 2:  $\mathbf{g}_0 = 0$ ;  $\mathbf{x}_0^* = \mathbf{x}$ ;
- 3: **for**  $t = 0$  to  $T - 1$  **do**
- 4:   Input  $\mathbf{x}_t^*$  to  $f$  and obtain the gradient  $\nabla_{\mathbf{x}} J(\mathbf{x}_t^*, y)$ ;
- 5:   Update  $\mathbf{g}_{t+1}$  by accumulating the velocity vector in the gradient direction as

$$\mathbf{g}_{t+1} = \mu \cdot \mathbf{g}_t + \frac{\nabla_{\mathbf{x}} J(\mathbf{x}_t^*, y)}{\|\nabla_{\mathbf{x}} J(\mathbf{x}_t^*, y)\|_1}; \quad (6)$$

- 6:   Update  $\mathbf{x}_{t+1}^*$  by applying the sign gradient as

$$\mathbf{x}_{t+1}^* = \mathbf{x}_t^* + \alpha \cdot \text{sign}(\mathbf{g}_{t+1}); \quad (7)$$

7: **end for**

8: **return**  $\mathbf{x}^* = \mathbf{x}_T^*$ .

---

<sup>20</sup>Dong, Yinpeng, et al. "Boosting adversarial attacks with momentum." 2017



# Сравнение FGSM-like атак

	Attack	Inc-v3	Inc-v4	IncRes-v2	Res-152	Inc-v3 <sub>ens3</sub>	Inc-v3 <sub>ens4</sub>	IncRes-v2 <sub>ens</sub>
Inc-v3	FGSM	72.3*	28.2	26.2	25.3	11.3	10.9	4.8
	I-FGSM	<b>100.0*</b>	22.8	19.9	16.2	7.5	6.4	4.1
	MI-FGSM	<b>100.0*</b>	<b>48.8</b>	<b>48.0</b>	<b>35.6</b>	<b>15.1</b>	<b>15.2</b>	<b>7.8</b>
Inc-v4	FGSM	32.7	61.0*	26.6	27.2	13.7	11.9	6.2
	I-FGSM	35.8	<b>99.9*</b>	24.7	19.3	7.8	6.8	4.9
	MI-FGSM	<b>65.6</b>	<b>99.9*</b>	<b>54.9</b>	<b>46.3</b>	<b>19.8</b>	<b>17.4</b>	<b>9.6</b>
IncRes-v2	FGSM	32.6	28.1	55.3*	25.8	13.1	12.1	7.5
	I-FGSM	37.8	20.8	<b>99.6*</b>	22.8	8.9	7.8	5.8
	MI-FGSM	<b>69.8</b>	<b>62.1</b>	99.5*	<b>50.6</b>	<b>26.1</b>	<b>20.9</b>	<b>15.7</b>
Res-152	FGSM	35.0	28.2	27.5	72.9*	14.6	13.2	7.5
	I-FGSM	26.7	22.7	21.2	<b>98.6*</b>	9.3	8.9	6.2
	MI-FGSM	<b>53.6</b>	<b>48.9</b>	<b>44.7</b>	98.5*	<b>22.1</b>	<b>21.7</b>	<b>12.9</b>



# Метод атаки: One pixel

- Однопиксельная атака<sup>21</sup> — предельный случай  $\ell_0$ -атаки

---

<sup>21</sup>Su, Jiawei, Danilo Vasconcellos Vargas, and Kouichi Sakurai. "One pixel attack for fooling deep neural networks." 2017

<sup>22</sup>Storn, Rainer, and Kenneth Price. "Differential evolution — a simple and efficient heuristic for global optimization over continuous spaces." 1997

# Метод атаки: One pixel

- Однопиксельная атака<sup>21</sup> — предельный случай  $\ell_0$ -атаки
- **Идея:** применить эволюционный алгоритм (дифференциальной эволюции<sup>22</sup>)

---

<sup>21</sup>Su, Jiawei, Danilo Vasconcellos Vargas, and Kouichi Sakurai. "One pixel attack for fooling deep neural networks." 2017

<sup>22</sup>Storn, Rainer, and Kenneth Price. "Differential evolution — a simple and efficient heuristic for global optimization over continuous spaces." 1997

# Метод атаки: One pixel

- Однопиксельная атака<sup>21</sup> — предельный случай  $\ell_0$ -атаки
- **Идея:** применить эволюционный алгоритм (дифференциальной эволюции<sup>22</sup>)
- Популяция состоит из 400 экземпляров, каждый из которых задается пятеркой: две координаты и три канала цвета

---

<sup>21</sup>Su, Jiawei, Danilo Vasconcellos Vargas, and Kouichi Sakurai. "One pixel attack for fooling deep neural networks." 2017

<sup>22</sup>Storn, Rainer, and Kenneth Price. "Differential evolution — a simple and efficient heuristic for global optimization over continuous spaces." 1997

# Метод атаки: One pixel

- Однопиксельная атака<sup>21</sup> — предельный случай  $\ell_0$ -атаки
- **Идея:** применить эволюционный алгоритм (дифференциальной эволюции<sup>22</sup>)
- Популяция состоит из 400 экземпляров, каждый из которых задается пятеркой: две координаты и три канала цвета
- Генерация потомка — линейная комбинация трех случайных родителей

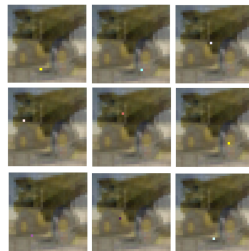
---

<sup>21</sup>Su, Jiawei, Danilo Vasconcellos Vargas, and Kouichi Sakurai. "One pixel attack for fooling deep neural networks." 2017

<sup>22</sup>Storn, Rainer, and Kenneth Price. "Differential evolution — a simple and efficient heuristic for global optimization over continuous spaces." 1997

# Метод атаки: One pixel

- Однопиксельная атака<sup>21</sup> — предельный случай  $\ell_0$ -атаки
- **Идея:** применить эволюционный алгоритм (дифференциальной эволюции<sup>22</sup>)
- Популяция состоит из 400 экземпляров, каждый из которых задается пятеркой: две координаты и три канала цвета
- Генерация потомка — линейная комбинация трех случайных родителей



Original Image (dog)

Airplane	Automobile	Bird
Cat	Deer	Frog
Horse	Ship	Truck

Target classes

<sup>21</sup>Su, Jiawei, Danilo Vasconcellos Vargas, and Kouichi Sakurai. "One pixel attack for fooling deep neural networks." 2017

<sup>22</sup>Storn, Rainer, and Kenneth Price. "Differential evolution — a simple and efficient heuristic for global optimization over continuous spaces." 1997



- Все атаки до этого работали в т.н. цифровой области (digital domain): изменяли картинку на уровне пикселей

---

<sup>23</sup>Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. "Adversarial examples in the physical world." 2016



- Все атаки до этого работали в т.н. цифровой области (digital domain): изменяли картинку на уровне пикселей
- Если нет возможности проатаковать изображение непосредственно перед подачей в СНС, то такая атака бесполезна

---

<sup>23</sup>Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. "Adversarial examples in the physical world." 2016



- Все атаки до этого работали в т.н. цифровой области (digital domain): изменяли картинку на уровне пикселей
- Если нет возможности проатаковать изображение непосредственно перед подачей в СНС, то такая атака бесполезна
- Поэтому атаки в реальном мире (real-world), или физические атаки, наиболее универсальны

---

<sup>23</sup>Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. "Adversarial examples in the physical world." 2016



- Все атаки до этого работали в т.н. цифровой области (digital domain): изменяли картинку на уровне пикселей
- Если нет возможности проатаковать изображение непосредственно перед подачей в СНС, то такая атака бесполезна
- Поэтому атаки в реальном мире (real-world), или физические атаки, наиболее универсальны
- Первый пример физической атаки<sup>23</sup> — атака на изображение в цифровой области, затем печать на физическом носителе (бумага), затем снимок цифровой камерой и последующая обработка СНС

---

<sup>23</sup>Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. "Adversarial examples in the physical world." 2016

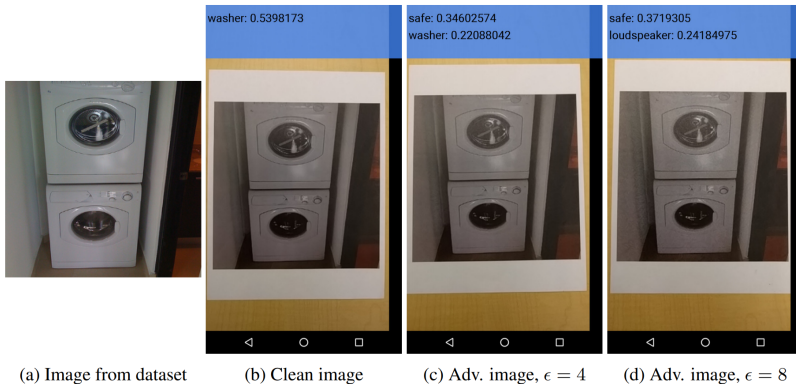


- Все атаки до этого работали в т.н. цифровой области (digital domain): изменяли картинку на уровне пикселей
- Если нет возможности проатаковать изображение непосредственно перед подачей в СНС, то такая атака бесполезна
- Поэтому атаки в реальном мире (real-world), или физические атаки, наиболее универсальны
- Первый пример физической атаки<sup>23</sup> — атака на изображение в цифровой области, затем печать на физическом носителе (бумага), затем снимок цифровой камерой и последующая обработка СНС
- Никакой специальной технологии для генерации таких атак еще не было, просто была показана их возможность

---

<sup>23</sup>Kurakin, Alexey, Ian Goodfellow, and Samy Bengio. "Adversarial examples in the physical world." 2016





- Подход EOT<sup>24</sup> (**E**xpectation **O**ver **T**ransformation) учитывает, что объект в реальном мире обычно претерпевает ряд преобразований таких как:
  - Масштабирование
  - Трансляция (тряска)
  - Изменение яркости и/или контрастности

---

<sup>24</sup>Athalye, Anish, et al. "Synthesizing robust adversarial examples." 2017



- Подход EOT<sup>24</sup> (Expectation Over Transformation) учитывает, что объект в реальном мире обычно претерпевает ряд преобразований таких как:
  - Масштабирование
  - Трансляция (тряска)
  - Изменение яркости и/или контрастности
- Поэтому задача — найти (направленную) состязательную атаку  $r$  с учетом множества преобразований  $T$ :

## EOT

Найти  $\arg \max_r \mathbb{E}_{g \sim T} P(y_t | g(x + r))$  при условии:

- 1  $f(x) = y_{gt} \neq y_t$
- 2  $\mathbb{E}_{g \sim T} \|g(x + r) - g(x)\|_p < \epsilon$
- 3  $x \in B$

<sup>24</sup>Athalye, Anish, et al. "Synthesizing robust adversarial examples." 2017

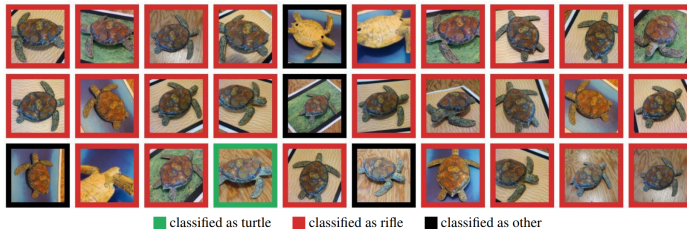


- В итоге, используя широкий ряд преобразований  $T$ , удалось сделать состязательный 3D-пример



- В итоге, используя широкий ряд преобразований  $T$ , удалось сделать состязательный 3D-пример

Transformation	Minimum	Maximum
Camera distance	2.5	3.0
X/Y translation	-0.05	0.05
Rotation	any	
Background	(0.1, 0.1, 0.1)	(1.0, 1.0, 1.0)
Lighten / Darken (additive)	-0.15	0.15
Lighten / Darken (multiplicative)	0.5	2.0
Per-channel (additive)	-0.15	0.15
Per-channel (multiplicative)	0.7	1.3
Gaussian Noise (stdev)	0.0	0.1



## Еще примеры физических атак

- Интересны примеры атак на объекты ImageNet<sup>25</sup>, дорожные знаки<sup>26</sup> и даже системы распознавания лиц<sup>27</sup>

---

<sup>25</sup>Brown, Tom B., et al. "Adversarial patch." 2017

<sup>26</sup>Eykholt, Kevin, et al. "Robust physical-world attacks on deep learning models." 2017

<sup>27</sup>Sharif, Mahmood, et al. "Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition." 2016



## Еще примеры физических атак

- Интересны примеры атак на объекты ImageNet<sup>25</sup>, дорожные знаки<sup>26</sup> и даже системы распознавания лиц<sup>27</sup>
- Примечательно, что все эти атаки по существу  $\ell_0$ -атаки, а также используют NPS и TV-добавки в функцию потерь
  - NPS (**N**on **P**rintability **S**core): штраф за использование цветов, которые не может воспроизвести данный принтер
  - TV (**T**otal **V**ariation): штраф за негладкость картинки

$$TV(x) = \sum_{i,j} \sqrt{(x_{i,j+1} - x_{i,j})^2 + (x_{i+1,j} - x_{i,j})^2}$$

<sup>25</sup>Brown, Tom B., et al. "Adversarial patch." 2017

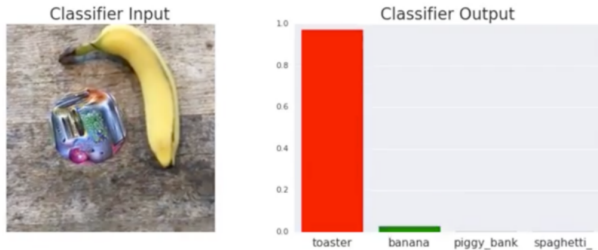
<sup>26</sup>Eykholt, Kevin, et al. "Robust physical-world attacks on deep learning models." 2017

<sup>27</sup>Sharif, Mahmood, et al. "Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition." 2016



# Еще примеры физических атак

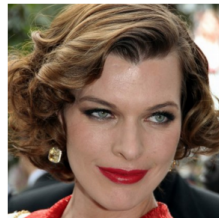
## Атака на объекты ImageNet:



## Атака на дорожные знаки:



## Атака на FaceID:



# Атаки на ведущую систему распознавания лиц

- Обычно система распознавания содержит два важных элемента: детектор и извлекатель признаков (часто называемый FaceID)

---

<sup>28</sup>Zhang, Kaipeng, et al. "Joint face detection and alignment using multitask cascaded convolutional networks." 2016

<sup>29</sup>Deng, Jiankang, et al. "Arcface: Additive angular margin loss for deep face recognition." 2018



# Атаки на ведущую систему распознавания лиц

- Обычно система распознавания содержит два важных элемента: детектор и извлекатель признаков (часто называемый FaceID)
- Использовались: крайне легкий нейросетевой детектор MTCNN<sup>28</sup> и ведущая открытая система извлечения признаков ArcFace<sup>29</sup>

---

<sup>28</sup>Zhang, Kaipeng, et al. "Joint face detection and alignment using multitask cascaded convolutional networks." 2016

<sup>29</sup>Deng, Jiankang, et al. "Arcface: Additive angular margin loss for deep face recognition." 2018



# Атаки на ведущую систему распознавания лиц

- Обычно система распознавания содержит два важных элемента: детектор и извлекатель признаков (часто называемый FaceID)
- Использовались: крайне легкий нейросетевой детектор MTCNN<sup>28</sup> и ведущая открытая система извлечения признаков ArcFace<sup>29</sup>
- Атаки на FaceID: с цветным патчем и черно-белым

---

<sup>28</sup>Zhang, Kaipeng, et al. "Joint face detection and alignment using multitask cascaded convolutional networks." 2016

<sup>29</sup>Deng, Jiankang, et al. "Arcface: Additive angular margin loss for deep face recognition." 2018





# Атаки на ведущую систему распознавания лиц

- Обычно система распознавания содержит два важных элемента: детектор и извлекатель признаков (часто называемый FaceID)
- Использовались: крайне легкий нейросетевой детектор MTCNN<sup>28</sup> и ведущая открытая система извлечения признаков ArcFace<sup>29</sup>
- Атаки на FaceID: с цветным патчем и черно-белым
- Атака на детектор: маска и черно-белый патч

---

<sup>28</sup>Zhang, Kaipeng, et al. "Joint face detection and alignment using multitask cascaded convolutional networks." 2016

<sup>29</sup>Deng, Jiankang, et al. "Arcface: Additive angular margin loss for deep face recognition." 2018



# Атака на детектор лиц

- МТСNN очень простой и неглубокий и поэтому крайне устойчивый к состязательным атакам детектор



# Атака на детектор лиц

- MTCNN очень простой и неглубокий и поэтому крайне устойчивый к состязательным атакам детектор
- В MTCNN каскадный подход: сначала грубое приближение (P-Net), а затем исправление (R-Net, O-Net)



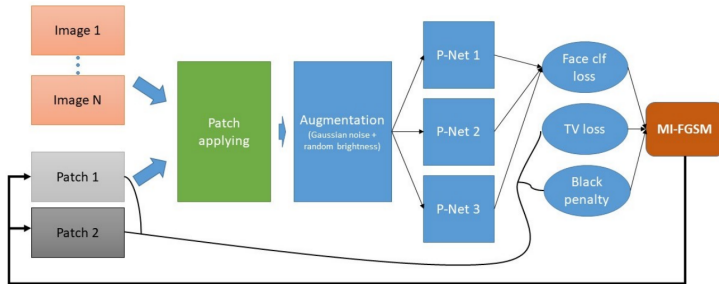
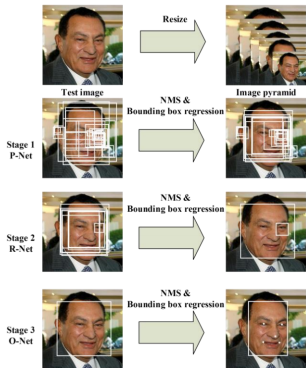
# Атака на детектор лиц

- MTCNN очень простой и неглубокий и поэтому крайне устойчивый к состязательным атакам детектор
- В MTCNN каскадный подход: сначала грубое приближение (P-Net), а затем исправление (R-Net, O-Net)
- Решение: атаковать самый первый и важный классификационный слой в P-Net



# Атака на детектор лиц

- MTCNN очень простой и неглубокий и поэтому крайне устойчивый к состязательным атакам детектор
- В MTCNN каскадный подход: сначала грубое приближение (P-Net), а затем исправление (R-Net, O-Net)
- Решение: атаковать самый первый и важный классификационный слой в P-Net



# Атака на детектор лиц<sup>31</sup> — черно-белые патчи

- Для физической атаки<sup>30</sup> пришлось оценивать параметры локальных проекций по заранее подготовленной маске

---

<sup>30</sup><https://www.youtube.com/watch?v=0Y700IS8bxs>

<sup>31</sup>Kaziakhmedov, Edgar, et al. "Real-world attack on MTCNN face detection system." 2019



# Атака на детектор лиц<sup>31</sup> — черно-белые патчи

- Для физической атаки<sup>30</sup> пришлось оценивать параметры локальных проекций по заранее подготовленной маске
- Из-за неглубокого характера детектора патчи не носят семантический характер

---

<sup>30</sup><https://www.youtube.com/watch?v=0Y700IS8bxs>

<sup>31</sup>Kaziakhmedov, Edgar, et al. "Real-world attack on MTCNN face detection system." 2019



# Атака на детектор лиц<sup>31</sup> — черно-белые патчи

- Для физической атаки<sup>30</sup> пришлось оценивать параметры локальных проекций по заранее подготовленной маске
- Из-за неглубокого характера детектора патчи не носят семантический характер



<sup>30</sup><https://www.youtube.com/watch?v=0Y700IS8bxs>

<sup>31</sup>Kaziakhmedov, Edgar, et al. "Real-world attack on MTCNN face detection system." 2019



- Т.н. “off-plane” проекция аналитически  $\Rightarrow$  можно пропускать градиенты

---

<sup>32</sup>Komkov, Stepan, and Aleksandr Petiushko. “AdvHat: Real-world adversarial attack on ArcFace Face ID system.” 2019

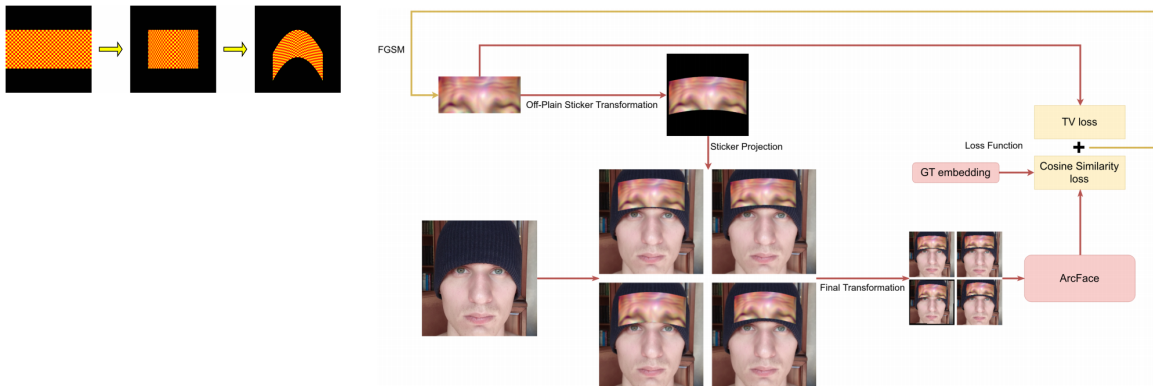
- Т.н. “off-plane” проекция аналитически  $\Rightarrow$  можно пропускать градиенты
- Патчи благодаря большой области восприятия носят семантический характер

---

<sup>32</sup>Komkov, Stepan, and Aleksandr Petiushko. “AdvHat: Real-world adversarial attack on ArcFace Face ID system.” 2019

# AdvHat<sup>32</sup> — шапка-невидимка

- Т.н. “off-plane” проекция аналитически  $\Rightarrow$  можно пропускать градиенты
- Патчи благодаря большой области восприятия носят семантический характер

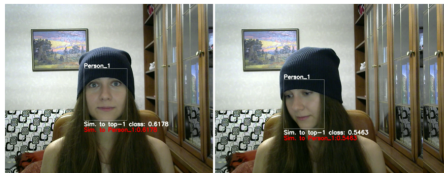


<sup>32</sup>Komkov, Stepan, and Aleksandr Petiushko. “AdvHat: Real-world adversarial attack on ArcFace Face ID system.” 2019

## Устойчивость к поворотам и разной освещенности<sup>33</sup>:

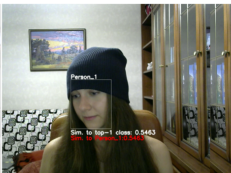
**Фронтальное лицо  
(нет атаки)**

Близость до своего эталона: **0.61**



**Поворот лица  
(нет атаки)**

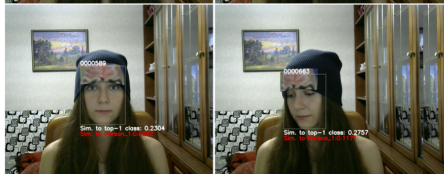
Близость до своего эталона: **0.54**



**Фронтальное лицо  
(атака)**

Близость до своего эталона: **0.02**

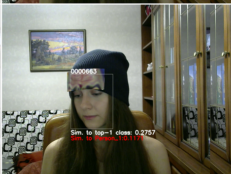
Близость до другого эталона: **0.23**



**Поворот лица  
(атака)**

Близость до своего эталона: **0.11**

Близость до другого эталона: **0.27**



<sup>33</sup><https://www.youtube.com/watch?v=a4iNg0wWBsQ>

# Adversarial patches<sup>34</sup> — черно-белые патчи

Дальнейшее развитие атак на FaceID:



<sup>34</sup>Pautov, Mikhail, et al. "On adversarial patches: real-world attack on ArcFace-100 face recognition system." 2019

# Вариант защиты от состязательных атак в реальном мире<sup>35</sup>

- Большинство состязательных атак в реальном мире основано на том, что к объекту добавляется специальная (обычно – прямоугольная) картинка, которая и ломает распознавание.

---

<sup>35</sup>Wu, Tong, Liang Tong, and Yevgeniy Vorobeychik. “Defending Against Physically Realizable Attacks on Image Classification.” 2019

# Вариант защиты от состязательных атак в реальном мире<sup>35</sup>

- Большинство состязательных атак в реальном мире основано на том, что к объекту добавляется специальная (обычно – прямоугольная) картинка, которая и ломает распознавание.
  - А давайте будем обучать в цифровой области (на обычных картинках), используя состязательное обучение и добавляя специальную прямоугольную аугментацию!

---

<sup>35</sup>Wu, Tong, Liang Tong, and Yevgeniy Vorobeychik. “Defending Against Physically Realizable Attacks on Image Classification.” 2019

# Вариант защиты от состязательных атак в реальном мире<sup>35</sup>

- Большинство состязательных атак в реальном мире основано на том, что к объекту добавляется специальная (обычно – прямоугольная) картинка, которая и ломает распознавание.
  - А давайте будем обучать в цифровой области (на обычных картинках), используя состязательное обучение и добавляя специальную прямоугольную аугментацию!
- Состязательное обучение предлагается делать двухэтапным:

---

<sup>35</sup>Wu, Tong, Liang Tong, and Yevgeniy Vorobeychik. “Defending Against Physically Realizable Attacks on Image Classification.” 2019



# Вариант защиты от состязательных атак в реальном мире<sup>35</sup>

- Большинство состязательных атак в реальном мире основано на том, что к объекту добавляется специальная (обычно – прямоугольная) картинка, которая и ломает распознавание.
  - А давайте будем обучать в цифровой области (на обычных картинках), используя состязательное обучение и добавляя специальную прямоугольную аугментацию!
- Состязательное обучение предлагается делать двухэтапным:
  - Сначала ищем наилучшую позицию для прямоугольника (среднего серого цвета),

<sup>35</sup>Wu, Tong, Liang Tong, and Yevgeniy Vorobeychik. “Defending Against Physically Realizable Attacks on Image Classification.” 2019

# Вариант защиты от состязательных атак в реальном мире<sup>35</sup>

- Большинство состязательных атак в реальном мире основано на том, что к объекту добавляется специальная (обычно – прямоугольная) картинка, которая и ломает распознавание.
  - А давайте будем обучать в цифровой области (на обычных картинках), используя состязательное обучение и добавляя специальную прямоугольную аугментацию!
- Состязательное обучение предлагается делать двухэтапным:
  - Сначала ищем наилучшую позицию для прямоугольника (среднего серого цвета),
    - Либо полным перебором (скользящим окном) по всем возможным позициям,

<sup>35</sup>Wu, Tong, Liang Tong, and Yevgeniy Vorobeychik. “Defending Against Physically Realizable Attacks on Image Classification.” 2019

# Вариант защиты от состязательных атак в реальном мире<sup>35</sup>

- Большинство состязательных атак в реальном мире основано на том, что к объекту добавляется специальная (обычно – прямоугольная) картинка, которая и ломает распознавание.
  - А давайте будем обучать в цифровой области (на обычных картинках), используя состязательное обучение и добавляя специальную прямоугольную аугментацию!
- Состязательное обучение предлагается делать двухэтапным:
  - Сначала ищем наилучшую позицию для прямоугольника (среднего серого цвета),
    - Либо полным перебором (скользящим окном) по всем возможным позициям,
    - Либо на основе позиций максимального значения градиента по входу,

<sup>35</sup>Wu, Tong, Liang Tong, and Yevgeniy Vorobeychik. “Defending Against Physically Realizable Attacks on Image Classification.” 2019

# Вариант защиты от состязательных атак в реальном мире<sup>35</sup>

- Большинство состязательных атак в реальном мире основано на том, что к объекту добавляется специальная (обычно – прямоугольная) картинка, которая и ломает распознавание.
  - А давайте будем обучать в цифровой области (на обычных картинках), используя состязательное обучение и добавляя специальную прямоугольную аугментацию!
- Состязательное обучение предлагается делать двухэтапным:
  - Сначала ищем наилучшую позицию для прямоугольника (среднего серого цвета),
    - Либо полным перебором (скользящим окном) по всем возможным позициям,
    - Либо на основе позиций максимального значения градиента по входу,
  - А затем – запускаем состязательную атаку (здесь PGD) внутри этого прямоугольника.



<sup>35</sup>Wu, Tong, Liang Tong, and Yevgeniy Vorobeychik. “Defending Against Physically Realizable Attacks on Image Classification.” 2019

# Заключительные выводы

- На данный момент СНС (в целом) работают гораздо лучше человека

---

<sup>36</sup>Image credit: <http://reddit.com>



# Заключительные выводы

- На данный момент СНС (в целом) работают гораздо лучше человека
- СНС легко “обмануть”, используя их неустойчивость по входу

---

<sup>36</sup>Image credit: <http://reddit.com>



# Заключительные выводы

- На данный момент СНС (в целом) работают гораздо лучше человека
- СНС легко “обмануть”, используя их неустойчивость по входу
- Наиболее распространенный прием атаки — производная по входу

---

<sup>36</sup>Image credit: <http://reddit.com>



# Заключительные выводы

- На данный момент СНС (в целом) работают гораздо лучше человека
- СНС легко “обмануть”, используя их неустойчивость по входу
- Наиболее распространенный прием атаки — производная по входу
- Перенести атаку в реальный мир непросто

---

<sup>36</sup>Image credit: <http://reddit.com>





# Заключительные выводы

- На данный момент СНС (в целом) работают гораздо лучше человека
- СНС легко “обмануть”, используя их неустойчивость по входу
- Наиболее распространенный прием атаки — производная по входу
- Перенести атаку в реальный мир непросто
- Однако можно сломать даже супер навороченные системы распознавания лиц, имея лишь обычный принтер

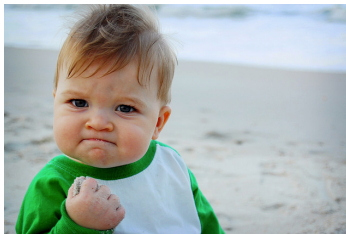
---

<sup>36</sup>Image credit: <http://reddit.com>



# Заключительные выводы

- На данный момент СНС (в целом) работают гораздо лучше человека
- СНС легко “обмануть”, используя их неустойчивость по входу
- Наиболее распространенный прием атаки — производная по входу
- Перенести атаку в реальный мир непросто
- Однако можно сломать даже супер навороченные системы распознавания лиц, имея лишь обычный принтер
- Для человечества пока еще не все потеряно<sup>36</sup>!



<sup>36</sup>Image credit: <http://reddit.com>

Спасибо за внимание!

