

Исследование устойчивости сверточных нейросетей на примере систем детекции и распознавания лиц

Петюшко Александр

МГУ им. М.В.Ломоносова, к.ф.-м.н.
Huawei, Video Intelligence Team Leader

4 февраля, 2021

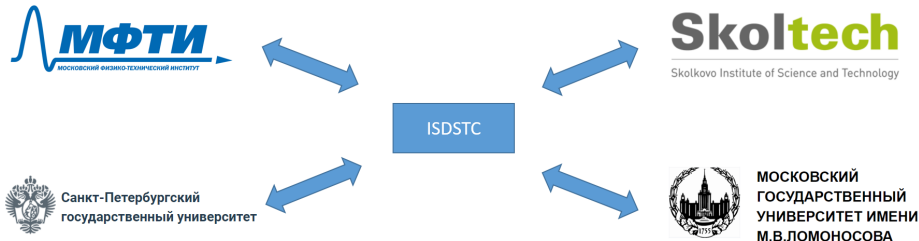


- 1 Лаборатория Интеллектуальных Систем и Науки о Данных
- 2 Потрясающие успехи СНС в компьютерном зрении
- 3 (Не) устойчивость СНС в компьютерном зрении
- 4 Методы для генерации состязательных примеров в цифровой области
- 5 ℓ_0 -состязательные примеры
- 6 Методы для генерации состязательных примеров в реальном мире
- 7 Состязательные примеры для систем детекции лиц
- 8 Состязательные примеры для систем распознавания лиц
- 9 Защита от состязательных примеров в реальном мире
- 10 Black-box восстановление лица по вектору признаков



Научное сотрудничество: Лаборатория Интеллектуальных Систем и Науки о Данных

Российский исследовательский институт → Московский исследовательский институт →
Лаборатория Интеллектуальных Систем и Науки о Данных



Человек или СНС?

ImageNet¹ (1000-классовая база данных изображений)

- Топ-5 ошибка для человека²: 5.1%
- Топ-5 ошибка для СНС³: 2.0%

Labeled Faces in the Wild⁴ (база данных лиц)

- Ошибка верификации для человека⁵: 2.47%
- Ошибка верификации для СНС⁶: 0.17%

¹<http://www.image-net.org/>

²<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/)

³Touvron, Hugo, et al. "Fixing the train-test resolution discrepancy." 2019

⁴<http://vis-www.cs.umass.edu/lfw/>

⁵Kumar, Neeraj, et al. "Attribute and simile classifiers for face verification." 2009

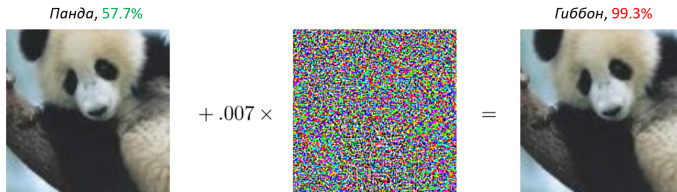
⁶Deng, Jiankang, et al. "Arcface: Additive angular margin loss for deep face recognition." 2018

LFW



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”⁷

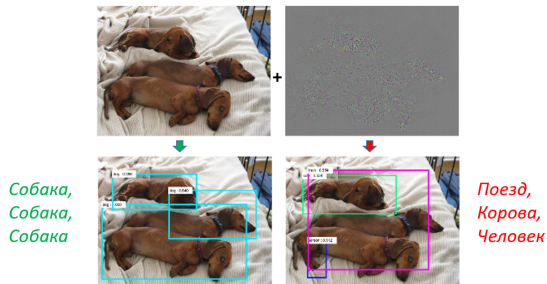


Такое возмущение называется **сопоставительным примером** (adversarial attack)

⁷Image credit: <https://arxiv.org/pdf/1412.6572.pdf>

Состязательные примеры в разных задачах

СНС для обнаружения и сегментации изображений⁸:



И даже НС для вопросно-ответных систем (QA, question answering systems)⁹:

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

⁸Xie C. et al. “Adversarial examples for semantic segmentation and object detection.” 2017

⁹Jia R. et al. “Adversarial examples for evaluating reading comprehension systems.” 2017

Состязательные примеры: 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 — правильный класс для x
- θ — параметры СНС-классификатора
- $L(\theta, x, y)$ — функция потерь
- $f(x)$ — выход классификатора (распознанный класс); при обучении мы добиваемся равенства $f(x) = y$
- $r \in B = [0, 1]^{C \times M \times N}$ — аддитивная добавка ко входу x



Состязательный пример и устойчивость: формулировки

Цель состязательного примера

Поменять выход классификатора f на неправильный путем добавления минимального (по некоторой норме ℓ_p) возмущения r :

- 1 $\|r\|_p \rightarrow \min$
- 2 $f(x) = y$ (изначальный ответ СНС верный)
- 3 $f(x + r) \neq y$ (меняем выход с помощью возмущения r)
- 4 $x + r \in B$ (остаемся во множестве допустимых изображений)

Устойчивость классификатора

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

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

Состязательные примеры в цифровой области: FGSM, итеративные методы

В основном состязательные примеры для СНС формулируются в терминах ℓ_∞ -нормы (что соответствует процессу восприятия человеческим глазом визуальной информации)

$$\|x\|_\infty = \max_i |x_i|, x = (x_1, \dots, x_n) \in \mathbb{R}^n$$

- **Fast Gradient Sign Method**¹⁰ (FGSM): $r = \epsilon \cdot \text{sign}(\nabla_x L(\theta, x, y))$
- Итеративный FGSM (I-FGSM)¹¹ / **Projected Gradient Descent** (PGD)¹² (где Π_B — операция проекции на B):

$$x^{t+1} = \Pi_B(x^t + \alpha \cdot \text{sign} \nabla_x L(\theta, x^t, y)), \quad x^0 = x, \alpha = \epsilon / T, t \leq T$$

- Сглаживание градиента для I-FGSM (MI-FGSM):

$$g^{t+1} = \mu \cdot g^t + \frac{\nabla_x L(\theta, x^t, y)}{\|\nabla_x L(\theta, x^t, y)\|_1}, x^{t+1} = \Pi_B(x^t + \alpha \cdot \text{sign}(g^{t+1})), \quad x^0 = x, g^0 = 0$$

¹⁰Goodfellow I. et al. "Explaining and harnessing adversarial examples." 2014

¹¹Kurakin A. et al. "Adversarial examples in the physical world." 2016

¹²Madry A. et al. "Towards deep learning models resistant to adversarial attacks." 2017



Сравнение состязательных примеров на основе FGSM

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

Как видно, MI-FGSM — наиболее успешная методика генерации.



ℓ_0 -состязательные примеры

- ℓ_∞ -состязательные примеры “заточены” под визуальное восприятие, но нужно обычно менять все пиксели
- В реальном мире это нереалистично — мы можем менять только часть сцены
- ℓ_0 -состязательные примеры более приспособлены для данной задачи: $\|x\|_0 = \sum_{i=1}^n \mathbf{1}_{x_i \neq 0}$
- **Jacobian-based Saliency Map Attack (JSMA)**¹³ и наиболее экстремальный случай — **One Pixel**¹⁴ — примеры подобных ℓ_0 -состязательных примеров, где минимизируется количество пикселей для изменения



¹³Papernot N. et al. “The limitations of deep learning in adversarial settings.” 2015

¹⁴Su J. et al. “One pixel attack for fooling deep neural networks.” 2017

Физические состязательные примеры: EOT

- Не имеем доступа к фото (и его пикселям) \Rightarrow единственная возможность — это изменить внешний вид самого объекта
- Подход **Expectation Over Transformation (EOT)**¹⁵ учитывает, что объект в реальном мире обычно претерпевает ряд преобразований таких как:
 - Масштабирование
 - Трансляции и повороты
 - Изменение яркости и/или контрастности, шум и т.п.
- Т.о. для объекта x нужно найти состязательный пример r с учетом преобразований $g \in T$:

EOT

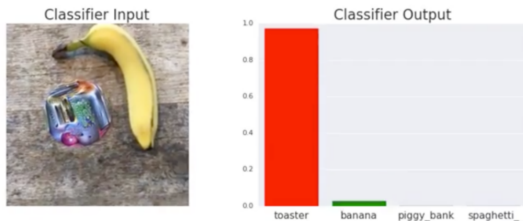
Найти $\arg \min_r \mathbb{E}_{g \sim T} [P(y|g(x+r))]$ при условии:

- 1 $\mathbb{E}_{g \sim T} [d(g(x+r), g(x))] < \epsilon$, где $d(a, b)$ – некоторая функция расстояние (например, $d(a, b) = \|a - b\|_p$)
- 2 $x + r \in B$

¹⁵Athalye A. et al. "Synthesizing robust adversarial examples." 2017

Физические состязательные примеры

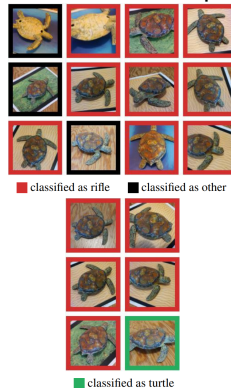
Для ImageNet¹⁶:



Для дорожных знаков¹⁷:



3D-состязательные примеры:



¹⁶Brown T. et al. "Adversarial patch." 2017

¹⁷Eykholt K. et al. "Robust physical-world attacks on deep learning models." 2017

Физические состязательные примеры: основные ингредиенты

- ℓ_0 оптимизация (маска) + EOT: обязательно
- **Total Variation (TV)** функция потерь — штраф за негладкость примера (в реальном мире мало больших попиксельных градиентов):

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

- **Non-Printability Score (NPS)** — штраф за использование цветов, которые не поддерживаются (принтером). Если $G \subset [0, 1]^3$ — поддерживаемая цветовая палитра (gamut), то штраф за использование цвета пикселя $q_0 \in [0, 1]^3$:

$$NPS(q_0) = \prod_{q \in G} \|q - q_0\|_2$$

- Дополнительный маппинг цветов (например, принтер печатает не цвет c , а близкий $m(c)$, и это нужно учитывать)



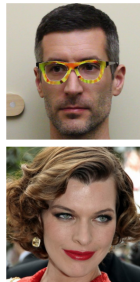
Состязательные примеры для систем детекции и распознавания лиц

- Изначально т.н. **Camouflage Art**¹⁸ использовался для обхода лидирующей системы детекции лиц Viola-Jones
- Это был вручную подобранный грим для обхода детектора Хаара
- Прорыв случился с работой Sharif et al.¹⁹, где предложили использовать состязательные очки
- Использовано: ℓ_0 -оптимизация + EOT + TV + NPS + цветовой маппинг
- Минусы: предобученные лица, системы распознавания лиц прошлого поколения

cvdazzle.com



Состязательные очки



¹⁸Feng R. et al. "Facilitating fashion camouflage art." 2013

¹⁹Sharif M. et al. "Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition." 2016



Общая схема детекции и распознавания лиц

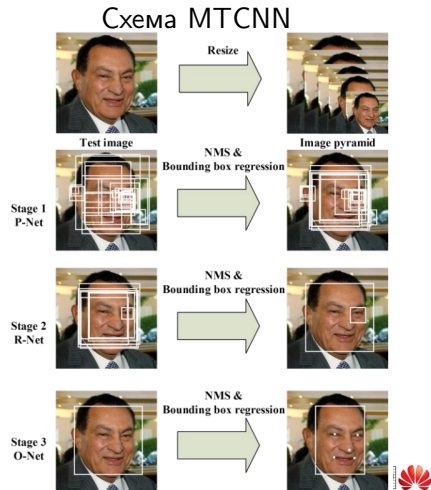


Общая схема детекции и распознавания лиц



Детектор лиц: MTCNN²⁰

- В отличие от современных глубоких детекторов типа Faster RCNN или YOLO, MTCNN очень простой и неглубокий \Rightarrow меньше поле восприятия, сложнее поменять решение детектора
- В MTCNN каскадный подход: сначала грубое приближение (P-Net), а затем исправление (R-Net, O-Net)
- Решение: использовать для генерации состязательных примеров самый первый классификационный слой в P-Net (а не функции потерь для прямоугольников или ключевых точек)

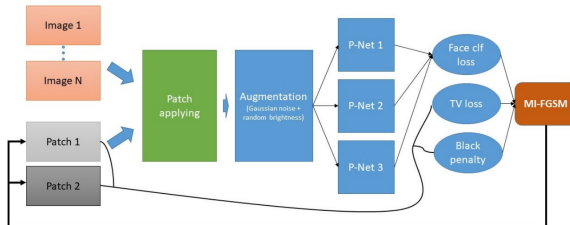


²⁰Zhang K. et al. "Joint face detection and alignment using multitask cascaded convolutional networks." 2016

Состязательные примеры для детектора лиц MTCNN

- EOT: Гауссов шум, размер маски, яркость, набор разных фото лица
- TV: +, NPS: –
- Маппинг цветов: штрафует за близость к черному цвету ($x_{i,j} = 1$) \Rightarrow новая добавка в функцию потерь: $L_{BLK}(x) = \sum_{i,j} x_{i,j}$
- Оптимизатор: MI-FGSM

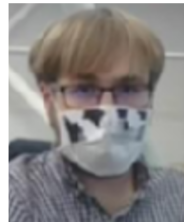
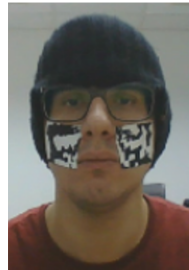
Схема генерации состязательных примеров для MTCNN



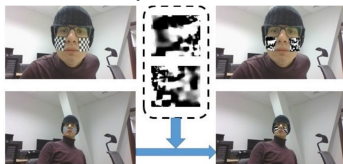
Состязательные примеры для детектора лиц MTCNN

- ℓ_0 -оптимизация: 2 версии
 - 1 две раздельных маски на щеках
 - 2 цельная медицинская маска
- У MTCNN маленькое поле восприятия \Rightarrow примеры не носят семантический характер
- Оценка параметров локальных аффинных проекций на основе заранее подготовленной специальной сетки для обучения

Примеры

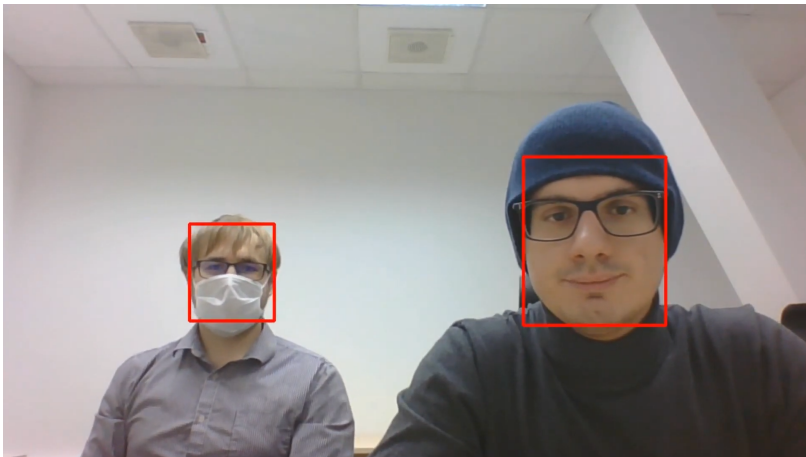


Проекция



Состязательные примеры для детектора лиц MTCNN: результат

Детали: статья²¹ (IEEE-2019) и видео-демонстрация²².



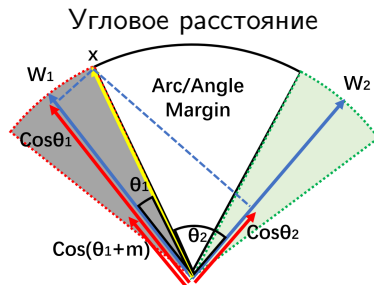
²¹Kaziakhmedov E. et al. "Real-world attack on MTCNN face detection system." 2019

²²<https://www.youtube.com/watch?v=0Y700IS8bxs>



Система распознавания лиц: ArcFace²³

- Выбрана ведущая открытая система распознавания лиц признаков: ArcFace
- Главная идея ArcFace — использовать угловое расстояние между векторами признаков
- Огромная обучающая база данных (MS1M) и глубокая СНС (ResNet-100)

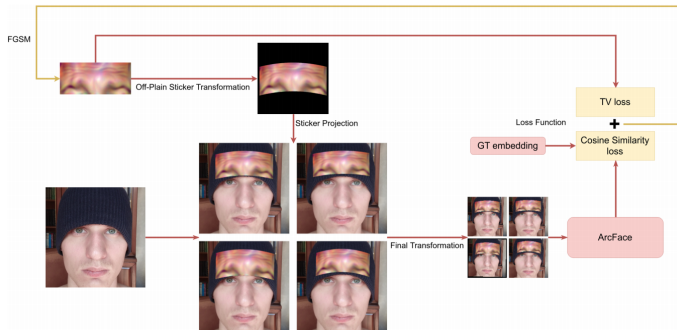


²³Deng J. et al. "Arcface: Additive angular margin loss for deep face recognition." 2018

Состязательные примеры для распознавателя лиц ArcFace

- EOT: различные параметры проекций маски, одно изображение лица
- TV: +, NPS: -, цветовой маппинг: -
- Добавка к функции потерь $L_{sim}(x, x_{gt}) = \cos(\text{emb}(x), \text{emb}(x_{gt}))$ для работы с любым лицом, где x_{gt} — фото лица, $\text{emb}(x)$ — вектор признаков x
- Оптимизатор: MI-FGSM

Схема генерации состязательных примеров для ArcFace



Состязательные примеры для распознавателя лиц ArcFace

- ℓ_0 -оптимизация: цветная связная область (патч) в области лба

- Глубокая СНС \Rightarrow большое поле восприятия \Rightarrow семантический характер примера

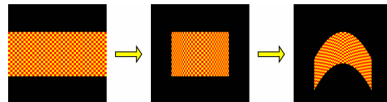
- Т.н. “off-plane” нелинейная проекция:

$$(x, y, 0) \rightarrow (x', y, z'), z' = a \cdot x'^2$$

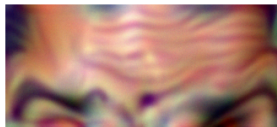
- Дифференцируемый слой Spatial Transformer Layer²⁴

$$x' = a \cdot \left(|x| \cdot \sqrt{x^2 + \frac{1}{4 \cdot a^2}} + \frac{1}{4 \cdot a^2} \cdot \ln \left(|x| + \sqrt{x^2 + \frac{1}{4 \cdot a^2}} \right) - \frac{1}{4 \cdot a^2} \cdot \ln \left(\frac{1}{2 \cdot a} \right) \right)$$

Off-plane проекция



Семантический характер примера



²⁴Jaderberg M. et al. “Spatial transformer networks.” 2015

Благодаря улучшенной процедуре проекции и полного использования цветовой палитры, состязательный пример устойчив к разным поворотам и освещенности

Анфас

(AdvHat: -)

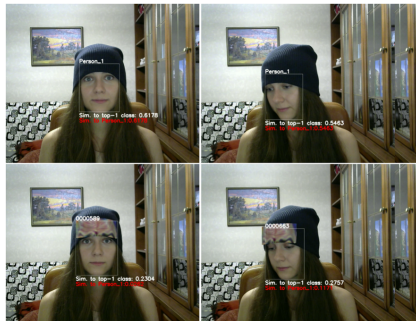
Близость (оригинал): **0.61**

Анфас

(AdvHat: +)

Близость (оригинал): **0.02**

Близость (другой): **0.23**



Профиль

(AdvHat: -)

Близость (оригинал): **0.54**

Профиль

(AdvHat: +)

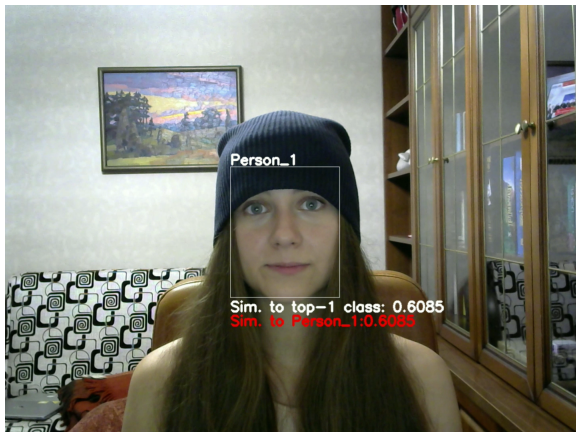
Близость (оригинал): **0.11**

Близость (другой): **0.27**



Состязательные примеры для распознавателя лиц ArcFace: результат

Детали: статья²⁵ (ICPR-2020) и видео-демонстрация²⁶.

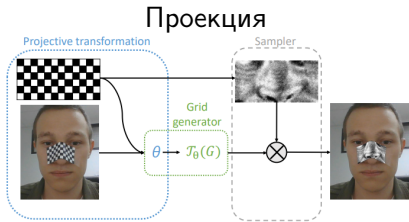


²⁵ Komkov S. et al. "Advhat: Real-world adversarial attack on arcface face id system." 2019

²⁶ <https://www.youtube.com/watch?v=a4iNg0wWBsQ>

Состязательные примеры для распознавателя лиц ArcFace: черно-белые патчи²⁷ (IEEE-2019)

- Комбинация двух предыдущих подходов:
 - Цветовой мappинг (штраф за черный цвет)
 - Локальные аффинные проекции на основе подготовленной сетки
- Пример носит семантический характер



Состязательные примеры



²⁷Pautov M. et al. "On adversarial patches: real-world attack on ArcFace-100 face recognition system." 2019

Защита от состязательных примеров для систем распознавания в реальном мире²⁹

- Большинство состязательных примеров в реальном мире — патчи
 - Предложение: **A**dversarial **T**raining (AT)²⁸ в пиксельной области + прямоугольная аугментация
- Поэтапное AT:
 - 1 Позиция серого патча (max функции потерь) via:
 - Полный перебор
 - Max градиента по входу
 - 2 MI-FGSM внутри этого патча для поиска раскраски

- Обычная процедура обучения:

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

- Состязательное обучение (AT):

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



²⁸Goodfellow I. et al. “Explaining and harnessing adversarial examples.” 2014

²⁹Wu T. et al. “Defending Against Physically Realizable Attacks on Image Classification.” 2019

Black-box восстановление лица по вектору признаков

- Black-box модель M : $M(x) = y$, где
 - x — входное фото лица
 - y — его векторное представление (эмбединг)
- **Задача:** восстановить x' , сохраняя личность (ID)
- **Ранее:** использование функции потерь MSE и метрик на признаках / GAN (NBNet³⁰)
- **Здесь:** метод оптимизации нулевого порядка для нахождения x' : $M(x') \approx M(x)$
 - В качестве M используется: ArcFace
 - Тест независимым критиком: FaceNET³¹
 - Функция потерь (близость):
 $1 - \cos(M(x'), M(x))$
 - Добавка: $(\|M(x')\|_2 - \|M(x)\|_2)^2$

- **Главная сложность:** огромное пространство поиска
- **Решение:** использование априорной информации о лице — 2D-гауссианы

$$G(x, y) = A \cdot e^{-\frac{(x-x_0)^2}{2\sigma_1^2} - \frac{(y-y_0)^2}{2\sigma_2^2}}$$



$$(x_0, y_0, \sigma_1, \sigma_2, A) = (56, 72, 22, 42, 1)$$

³⁰Mai G. et al. "On the reconstruction of face images from deep face templates." 2018

³¹Schroff F. et al. "Facenet: A unified embedding for face recognition and clustering." 2015.

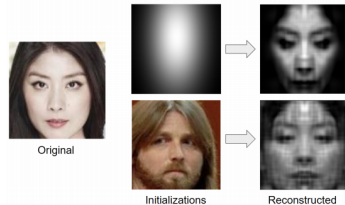
Black-box восстановление лица по вектору признаков: трюки

- Даже использование априорной информации недостаточно
- **Трюк1:** Использовать вертикальную симметрию лица \Rightarrow ищем только половину лица
- **Трюк2:** Для сохранения личности обычно цвет не нужен \Rightarrow ищем только 1 цветовой канал вместо 3



Original ArcFace: 0.978 ArcFace: 0.992 ArcFace: 0.961
FaceNet: 0.721 FaceNet: 0.685 FaceNet: 0.314

- **Инициализация:** Что взять за начальное приближение?
- **Обычно:** другое лицо (решение может быть смещено)
- **Можно:** оптимальный 2D-гауссиан (как раз нужна добавка на разность норм)



Симметрия, без симметрии, в цвете

Black-box восстановление лица по вектору признаков: алгоритм и примеры

Алгоритм




















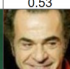
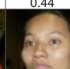

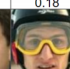
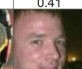
Algorithm 1 Face recovery algorithm

INPUT: target face embedding y , black-box model M , loss function L , $N_{queries}$

```
1:  $X \leftarrow 0$ 
2: Initialize  $G_0$ 
3: for  $i \leftarrow 0$  to  $N_{queries}$  do:
4:   Allocate image batch  $\mathbf{X}$ 
5:   Sample batch  $\mathbf{G}$  of random gaussians
6:    $\mathbf{X}_j = X + G_0 + \mathbf{G}_j$ 
7:    $\mathbf{y}' = M(\mathbf{X})$ 
8:    $\text{ind} = \text{argmin} \left( L(\mathbf{y}', y) \right)$ 
9:    $X \leftarrow X + \mathbf{G}_{\text{ind}}$ 
10:   $G_0 \leftarrow 0.99 \cdot G_0$ 
11:   $i \leftarrow i + \text{batchsize}$ 
12: end for
13:  $X \leftarrow X + G_0$ 
```

OUTPUT: reconstructed face X

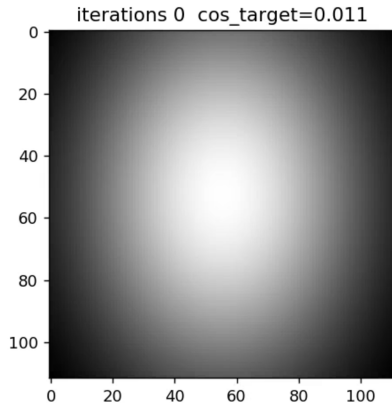
Примеры

Our method:						
ArcFace:	0.97	0.97	0.94	0.97	0.85	0.73
FaceNet:	0.70	0.75	0.72	0.78	0.38	-0.09
NBNet (WB):						
ArcFace:	0.17	0.21	0.12	0.26	0.06	0.09
FaceNet:	0.02	0.32	0.25	0.46	-0.01	0.35
NBNet (RGB):						
ArcFace:	0.28	0.46	0.34	0.54	0.12	0.21
FaceNet:	0.59	0.53	0.44	0.74	0.18	0.41
Original:						



Black-box восстановление лица по вектору признаков: результат

Детали: статья³² (ECCV-2020) и видео-демонстрация³³.



³²Razzhigaev A. et al. "Black-Box Face Recovery from Identity Features." 2020

³³<https://www.youtube.com/watch?v=sOrTcqRTw2A>

Заключительные выводы

- На данный момент СНС (в целом) работают гораздо лучше человека
- СНС неустойчивы (по входу)
- Перенести состязательный пример в реальный мир непросто
- Однако можно сломать даже супер навороченные системы распознавания лиц, имея лишь черно-белую бумажку с обычного принтера
- ℓ_0 -оптимизация + EOT + TV обязательны
- Через проекции нужно уметь пропускать градиенты
- Состязательное обучение может помочь в защите от состязательных примеров
- Изображения лиц могут быть восстановлены по векторам признаков даже в black-box манере



Присоединяйтесь к нам!

Кого мы ждем:

- ❑ Выпускники аспирантуры 2019-2021 годов
- ❑ Победители и призеры таких международных соревнований как ICPC, IMC, CTF, Kaggle, IMO, IOI, ICHO, IPHO etc.
- ❑ Техническое образование (информационные технологии, математика, физика, радиотехника, системы связи, информационная безопасность и др.)
- ❑ Английский на уровне "intermediate" и выше



Москва



Санкт-Петербург



Нижний Новгород



Новосибирск

Наши направления:

- Nonlinear algorithm development
- Wireless communication technologies
- Computer Vision with Deep Learning
- Math Library optimization
- Automatic program repair
- Compiler optimizations
- Automatic speech recognition
- AI databases and AI enabled systems
- Distributed and Parallel software
- Image/Video signal processing
- Software engineering and innovation
- Automated machine learning & Model optimization
- Computer architecture

Резюме можно выслать на почту:

rrhr@huawei.com



<https://career.huawei.ru/rri/>

Inspired by science to connect the world!

Looking forward to seeing your application



Спасибо за внимание!

