



НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ
УНИВЕРСИТЕТ

DeepFool: a simple and accurate method to fool deep neural networks

Арсения Шихова

Национальный исследовательский университет
«Высшая школа экономики»

13 декабря 2018г.

- Нейронная сеть обучается на конечном наборе изображений
- Злоумышленник знает всё: входные данные, архитектуру, параметры
- Цель: немного исказить картинку из трейна так, чтобы получить другой класс

- Беспилотные автомобили
- Всевозможные роботы и дроны
- Распознавание речи (голосовых команд)
- И так далее

- Нейронная сеть — это функция $f(x) = y; f : \mathbb{R}^n \rightarrow \mathbb{R}^c$, c — кол-во классов
- $\hat{k}(x) = \arg \max_{1 \leq i \leq c} y_i$ — класс
- Цель — $r : \hat{k}(x + r) \neq \hat{k}(x)$ — состязательный пример (adversarial example)

- Нейронная сеть — это функция $f(x) = y; f : \mathbb{R}^n \rightarrow \mathbb{R}^c$, c — кол-во классов
- $\hat{k}(x) = \arg \max_{1 \leq i \leq c} y_i$ — класс
- Цель — $r : \hat{k}(x + r) \neq \hat{k}(x)$ — состязательный пример (adversarial example)
- $\triangle(x, \hat{k}) = \min_r ||r||_2 : \hat{k}(x + r) \neq \hat{k}(x)$ — устойчивость (robustness) \hat{k} в точке x
- $\rho_{adv}(\hat{k}) = \mathbb{E}_x \frac{\triangle(x, \hat{k})}{||x||_2}$ — устойчивость классификатора \hat{k}

- Пусть $\hat{k}(x) = \text{sign}(f(x))$, где $f(x) = w^T x + b$
- Обозначим $\mathcal{F} = \{x : f(x) = 0\}$ — гиперплоскость
- Тогда $\Delta(x_0, \hat{k})$ — расстояние между x_0 и \mathcal{F}
- $r(x_0) = \arg \min_{r: \hat{k}(x_0+r) \neq \hat{k}(x_0)} \|r\|_2 = -\frac{f(x_0)}{\|w\|_2} w$
- $x_0 + r(x_0) = \frac{-b}{\|w\|_2^2} w$



Figure 2: Adversarial examples for a linear binary classifier.



- Пусть $\hat{k}(x) = \text{sign}(f(x))$
- Итеративно повышаем устойчивость $\Delta(x_0, \hat{k})$
- На i -м шаге $x_i = x_{i-1} + r_i$ с минимальным $\|r_i\|_2$, то есть $f(x_i) + \nabla f(x_i)^T r_i = 0$
- $r = \sum_i r_i$
- Ровно на границе быть плохо, в конце домножаем r на $1 + \eta$

Алгоритм 1: Binary DeepFool**Вход:** $x \in \mathbb{R}^n$, $\hat{k}(x) = \text{sign}(f(x))$ **Выход:** r : $\hat{k}(x) \neq \hat{k}(x + r)$ **Function** *DeepFool*(x, f): $x_0 = x$ $i = 0$ **while** $\text{sign}(f(x_i)) = \text{sign}(f(x_0))$ **do** $r_i = -\frac{f(x_i)}{\|\nabla f(x_i)\|_2^2} \nabla f(x_i)$ $x_{i+1} = x_i + r_i$ $i = i + 1$ **end****return** $r = \sum_i r_i$ **end**

- Используем подход **one-vs-all**
- $f : \mathbb{R}^n \rightarrow \mathbb{R}^c$
- $\hat{k}(x) = \arg \max_{1 \leq k \leq c} f_k(x)$
- $f(x) = W^T x + b$

- Используем подход **one-vs-all**
- $f : \mathbb{R}^n \rightarrow \mathbb{R}^c$
- $\hat{k}(x) = \arg \max_{1 \leq k \leq c} f_k(x)$
- $f(x) = W^T x + b$
- Цель:
$$\arg \min_{r: \exists k \neq \hat{k}(x_0) \ w_k^T(x_0+r) + b_k \geq w_{\hat{k}(x_0)}^T(x_0+r) + b_{\hat{k}(x_0)}} \|r\|_2^2$$
- Геометрическое объяснение: ищем расстояние между границей полиэдра P и точкой x_0
- $$P = \bigcap_{k=1}^c \{x : f_{\hat{k}(x_0)}(x) \geq f_k(x)\}$$
- Номер ближайшей к x_0 грани P есть $\hat{l}(x_0) = \arg \min_{k \neq \hat{k}(x_0)} \frac{|f_k(x_0) - f_{\hat{k}(x_0)}(x_0)|}{\|w_k - w_{\hat{k}(x_0)}\|_2}$
- $$r(x_0) = \frac{|f_{\hat{l}(x_0)}(x_0) - f_{\hat{k}(x_0)}(x_0)|}{\|w_{\hat{l}(x_0)} - w_{\hat{k}(x_0)}\|_2^2} (w_{\hat{l}(x_0)} - w_{\hat{k}(x_0)})$$
 — длина проекции x_0 на ближайшую грань P

- $f : \mathbb{R}^n \rightarrow \mathbb{R}^c$
- $\hat{k}(x) = \arg \max_{1 \leq k \leq c} f_k(x)$
- $P = \bigcap_{k=1}^c \{x : f_{\hat{k}(x_0)}(x) \geq f_k(x)\}$
- На каждой итерации минимизируем расстояние до
$$P_i = \bigcap_{k=1}^c \{x : f_k(x_i) - f_{\hat{k}(x_0)}(x_i) + \nabla f_k(x_i)^T x - \nabla f_{\hat{k}(x_0)}(x_i)^T x \leq 0\}$$

Алгоритм 2: Multiclass DeepFool

Вход: $x \in \mathbb{R}^n$, $\hat{k}(x) = \text{sign}(f(x))$

Выход: r : $\hat{k}(x) \neq \hat{k}(x + r)$

Function *DeepFool*(x, f):

```

     $x_0 = x$ 
     $i = 0$ 
    while  $\hat{k}(x_i) \neq \hat{k}(x_0)$  do
        for  $k \neq \hat{k}(x_0)$  do
             $w'_k = \nabla f_k(x_i) - \nabla f_{\hat{k}(x_0)}(x_i)$ 
             $f'_k = f_k(x_i) - f_{\hat{k}(x_0)}(x_i)$ 
        end
         $\hat{j} = \arg \min_{k \neq \hat{k}(x_0)} \frac{|f'_k|}{\|w'_k\|_2}$ 
         $r_i = \frac{|f'_j|}{\|w'_j\|_2^2} w'_j$ 
         $x_{i+1} = x_i + r_i$ 
         $i = i + 1$ 
    end
    return  $r = \sum_i r_i$ 
end

```

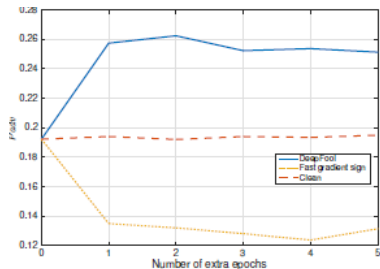
- Другой метод поиска примеров
- $r = \varepsilon \cdot \text{sign}(J(x, y, \theta))$, где θ — веса модели, J — стоимость обучения
- Подбираем ε так, чтобы исказить хотя бы 90% обучающей выборки

Classifier	Test error	$\hat{\rho}_{adv}$ [DeepFool]	time	$\hat{\rho}_{adv}$ [4]	time	$\hat{\rho}_{adv}$ [18]	time
LeNet (MNIST)	1%	2.0×10^{-1}	110 ms	1.0	20 ms	2.5×10^{-1}	> 4 s
FC500-150-10 (MNIST)	1.7%	1.1×10^{-1}	50 ms	3.9×10^{-1}	10 ms	1.2×10^{-1}	> 2 s
NIN (CIFAR-10)	11.5%	2.3×10^{-2}	1100 ms	1.2×10^{-1}	180 ms	2.4×10^{-2}	> 50 s
LeNet (CIFAR-10)	22.6%	3.0×10^{-2}	220 ms	1.3×10^{-1}	50 ms	3.9×10^{-2}	> 7 s
CaffeNet (ILSVRC2012)	42.6%	2.7×10^{-3}	510 ms*	3.5×10^{-2}	50 ms*	-	-
GoogLeNet (ILSVRC2012)	31.3%	1.9×10^{-3}	800 ms*	4.7×10^{-2}	80 ms*	-	-

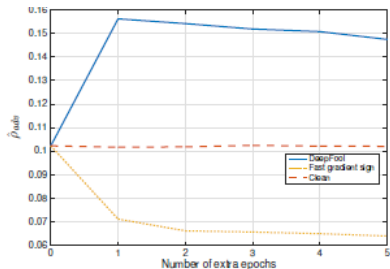
Table 1: The adversarial robustness of different classifiers on different datasets. The time required to compute one sample for each method is given in the time columns. The times are computed on a Mid-2015 MacBook Pro without CUDA support. The asterisk marks determines the values computed using a GTX 750 Ti GPU.

$$\text{здесь } \hat{\rho}_{adv}(f) = \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \frac{\|f(x)\|_2}{\|x\|_2}$$

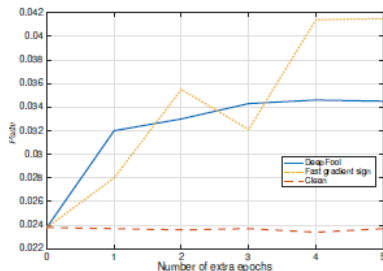
- Добавляем при обучении сгенерированные состязательные примеры
- Увеличиваем время обучения



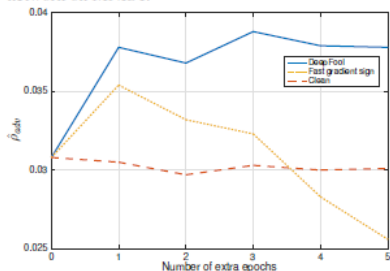
a) Effect of fine-tuning on adversarial examples computed by two different methods for LeNet on MNIST.



(b) Effect of fine-tuning on adversarial examples computed by two different methods for a fully-connected network on MNIST.



c) Effect of fine-tuning on adversarial examples computed by two different methods for NIN on CIFAR-10.



(d) Effect of fine-tuning on adversarial examples computed by two different methods for LeNet on CIFAR-10.

- cv-foundation.org/openaccess/content_cvpr_2016/papers/Moosavi-Dezfooli_DeepFool_A_Simple_CVPR_2016_paper.pdf — оригинал статьи
- arxiv.org/pdf/1412.6572.pdf — другой метод поиска состязательных примеров
- github.com/lts4/deepfool — репозиторий с кодом алгоритма