

Vulnerabilities of Deep Neural Networks

Alexey Gruzdev
a.s.gruzdev@yandex.ru

HSE, Summer, 2017

Can we fool the Neural Networks?



Short Answer:
Yes, we can

Long Answer: This work

Intriguing properties of neural networks - Feb 2014



x

“panda”

57.7% confidence

+ .007 ×



$\text{sign}(\nabla_x J(\theta, x, y))$

“nematode”

8.2% confidence

=



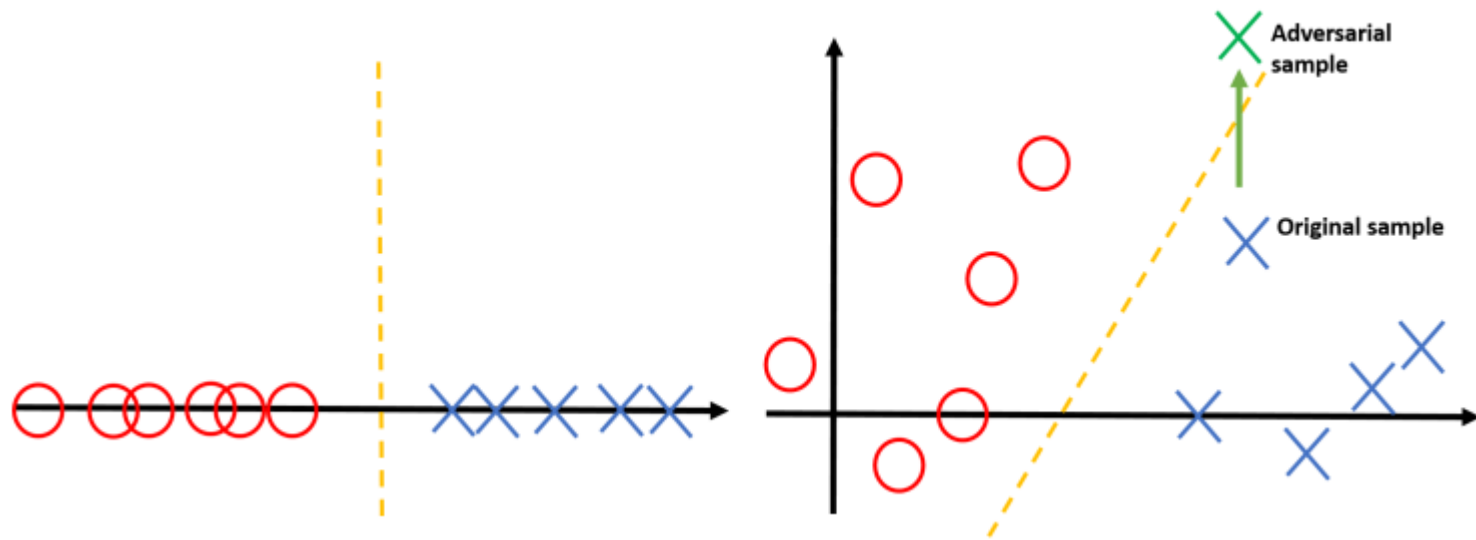
$x +$

$\epsilon \text{sign}(\nabla_x J(\theta, x, y))$

“gibbon”

99.3 % confidence

Adversarial intuition



Ok. But How?

- Fast sign (gradient sign) method:

$$\mathbf{X}^{adv} = \mathbf{X} + \epsilon \operatorname{sign}(\nabla_{\mathbf{X}} J(\mathbf{X}, y_{true}))$$

- Basic iterative method:

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

- Iterative least-likely method:

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

Clipping procedure

$Clip_{X,\epsilon} \{X'\}$ - function which performs per-pixel clipping of the image X' , so the result will be in L_∞ ϵ -neighbourhood of the source image X . The exact clipping equation is as follows:

$$Clip_{X,\epsilon} \{X'\} (x, y, z) = \min \left\{ 255, X(x, y, z) + \epsilon, \max \{ 0, X(x, y, z) - \epsilon, X'(x, y, z) \} \right\}$$

where $X(x, y, z)$ is the value of channel z of the image X at coordinates (x, y) .

Let's observe FGSM practically

original



eps=1.0



eps=2.0



eps=3.0



eps=10.0



Let's observe FGSM practically

original



eps=1.0



eps=2.0



eps=3.0



eps=10.0



Adversarial examples vs Original (FGSM)



Wait, we know all weights in DNN, what if we don't?

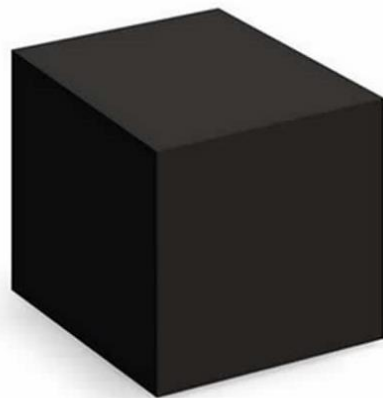
Q: What if we don't have access to model weights/ architecture & training data?

A: Let's approximate this model by another neural network and construct adversarial examples for `substitution`!

Black-Box attack

Black-Box attack:

1. Train substitute model (maybe not a neural network) on data with labels from oracle model
2. Construct adversarial examples for substitute model
3. Apply these examples for targeted model.



Results for ImageNet dataset

Attack Method	VGG-16	AlexNet	ResNet-50	Inception v3
FGSM	0.85/0.56/0.61	0.56/0.22/0.33	0.88/0.53/0.62	0.92/0.01/0.07
Black-Box	0.85/0.61	0.56/0.43	0.88/0.69	0.92/0.19

- FGSM is better for generating adversarial examples than Black-Box
- Inception v3 - best performer is much more vulnerable
- No architectures were stable

Results for CIFAR-10, CIFAR-100

Attack Method	VGG-16	ResNet-50
FGSM	0.79/0.03/0.36	0.81 /0.05/0.43
Black-Box	0.79/0.44	0.81 /0.51

Attack Method	VGG-16	ResNet-50
FGSM	0.64/0.15/0.36	0.67 /0.29/0.54
Black-Box	0.64/0.41	0.67 /0.38

Motivation:

- For ImageNet task input resolution is 224x224, what if you have less pixels to manipulate
- We still were able to drop the accuracy!

Transferability check

Attack Method	AlexNet	ResNet-50	Inception v3
FGSM	0.56/0.10/0.23	0.88/0.57/0.58	0.92/0.03/0.06

Attack Method	VGG-16	ResNet-50	Inception v3
FGSM	0.85/0.4/0.49	0.88/0.45/0.46	0.92/0.04/0.07

- Adversarial examples preserve the transferability - images generated to hack 1 model can hack others as well
- What's wrong with InceptionV3?!

Vulnerabilities of Deep Neural Networks

Alexey Gruzdev
a.s.gruzdev@yandex.ru

HSE, Summer, 2017