

Deep Learning - Adversarial attacks and defenses

Barbara Hammer

Machine Learning Group, Bielefeld University, Deep Learning 2024

Literature

- Adversarial Attacks and Defences: A Survey, Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, Debdeep Mukhopadhyay, 2018, <https://arxiv.org/abs/1810.00069>
- A Survey On Universal Adversarial Attack, Chaoning Zhang, Philipp Benz, Chenguo Lin, Adil Karjauv, Jing Wu, In So Kweon, 2021, <https://arxiv.org/abs/2103.01498>
- CAAI Transactions on Intelligence Technology, Volume 6, Issue 1 p. 25-45, A survey on adversarial attacks and defences, Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, Debdeep Mukhopadhyay, 2021 <https://doi.org/10.1049/cit2.12028>

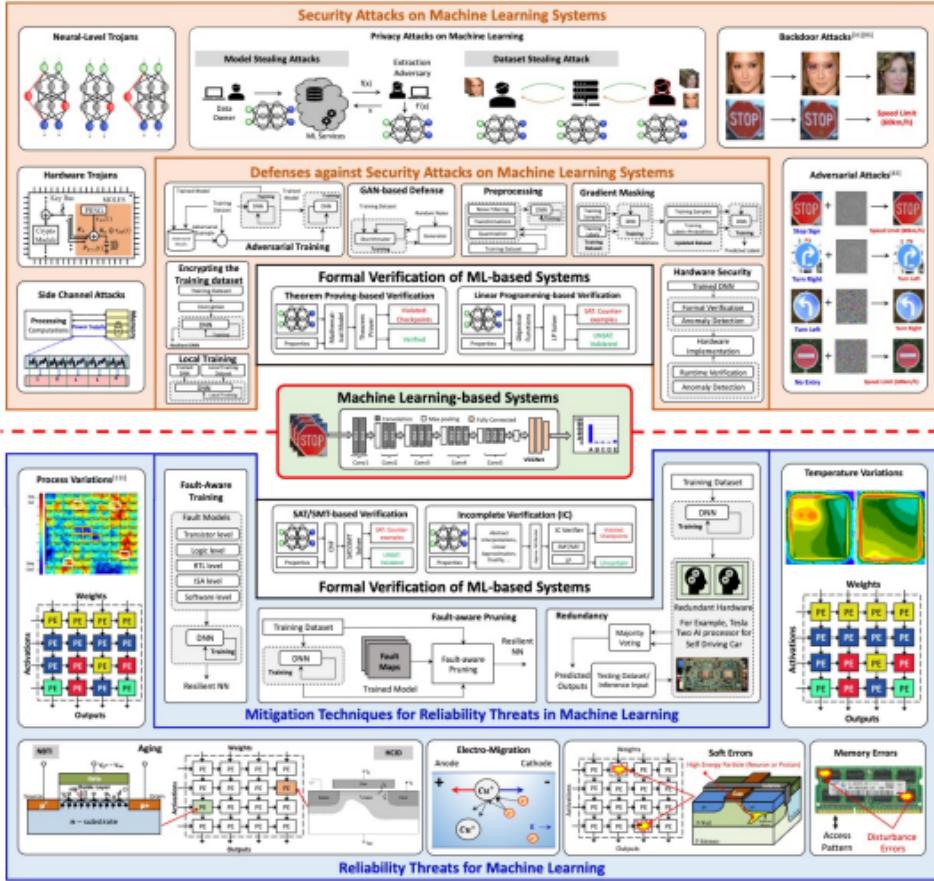
Code suite

- <https://github.com/cleverhans-lab/cleverhans>
- <https://deeprobust.readthedocs.io/en/latest/index.html>
- <https://adversarial-robustness-toolbox.readthedocs.io/en/latest/>

Content

- Survey of different types of attacks
- Adversarial attacks
 - Definitions and types
 - Adversarial attacks in practice
 - What makes an attack adversarial (rather than an error)
- Defenses against adversarial attacks
 - Defenses against adversarial attacks
 - Robust learning
 - Provably robust models

Attacks versus errors



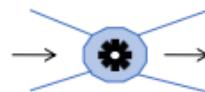
M. Shafique et al., "Robust Machine Learning Systems: Challenges, Current Trends, Perspectives, and the Road Ahead," in *IEEE Design & Test*, vol. 37, no. 2, pp. 30-57, April 2020, doi: 10.1109/MDAT.2020.2971217.

Types of attacks

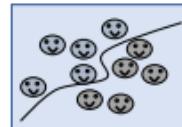
Learning pipeline



training data



learning algorithm

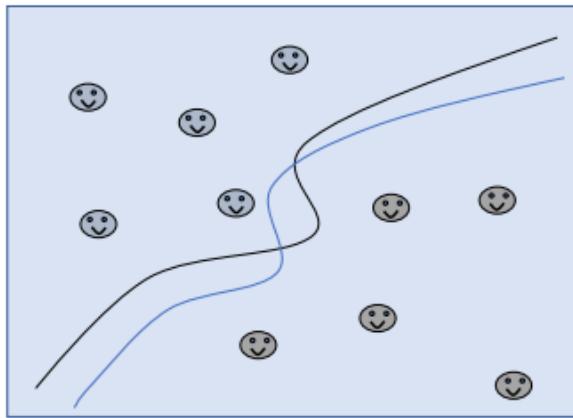


neural network



- speech processing
- computer vision
- machine translation
- ...

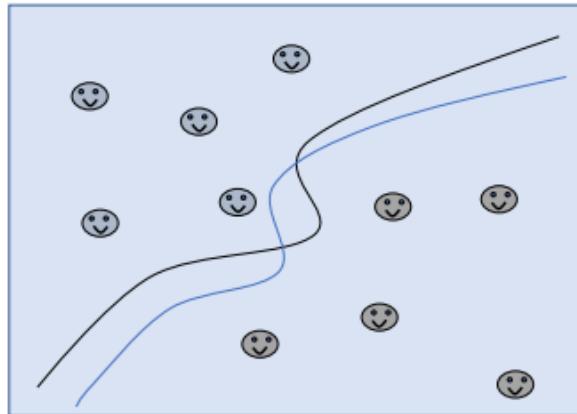
Learning theory



- train on representative data + guarantee smoothness of model → model behaves well for future data which are similar to the training data with high probability – **generalization ability of neural networks**

... but

- model depends on data which should be representative
- model is not exactly identical to the ground truth, i.e. it is always wrong at some places
- model is usually not causal but based on correlations



Where to attack?



Training a model

Using a model

Where to attack?



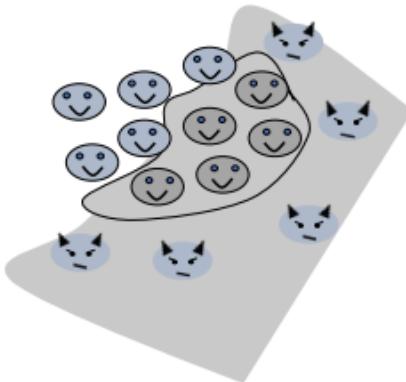
Training a model

- **data poisoning:** add training data such that the outcome is changed
- i.e. training data are not representative for what should be learned

Data Poisoning



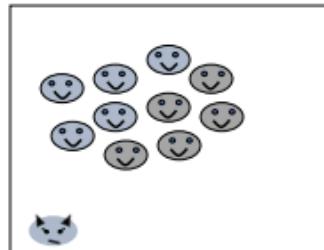
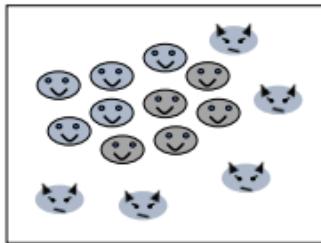
- **data poisoning:** add training data such that the outcome is changed



Data Poisoning



- **targeting model availability:** add malicious examples such that the trained model becomes useless – *user can observe poisoning*
- **targeting model integrity:** add malicious examples such that functional change stays hidden until used – *backdoor attacks*



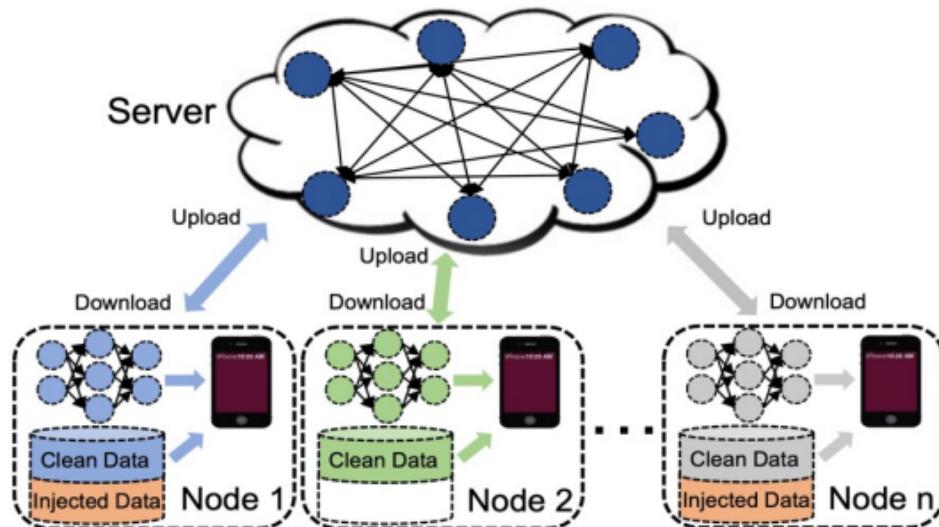
Making models useless

- **Tay:** AI chatter bot released by Microsoft via Twitter in 2016 to learn from interactions with humans, shut down after 16 hours due to offensive tweets
...



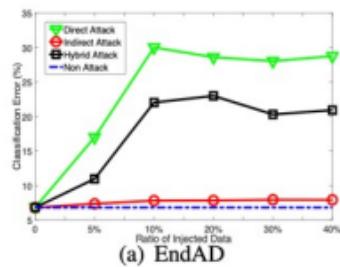
Twitter profile picture of Tay

Poisoning attack on federated learning

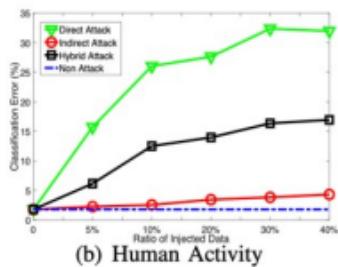


Data Poisoning Attacks on Federated Machine Learning, Gan Sun, Yang Cong (Senior Member, IEEE), Jiahua Dong, Qiang Wang, Ji Liu, 2021, <https://arxiv.org/abs/2004.10020>

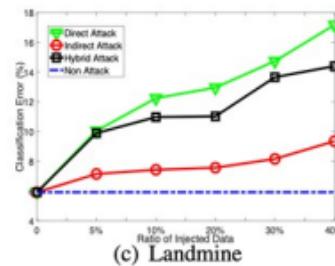
Poisoning attack on federated learning



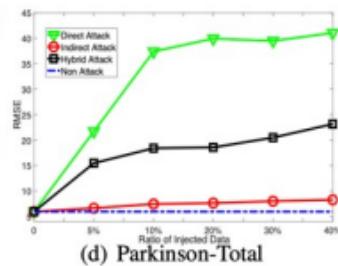
(a) EndAD



(b) Human Activity



(c) Landmine



(d) Parkinson-Total

direct attack:

introduce poisoned data to all target nodes in the network
(where performance is measured)

indirect attack:

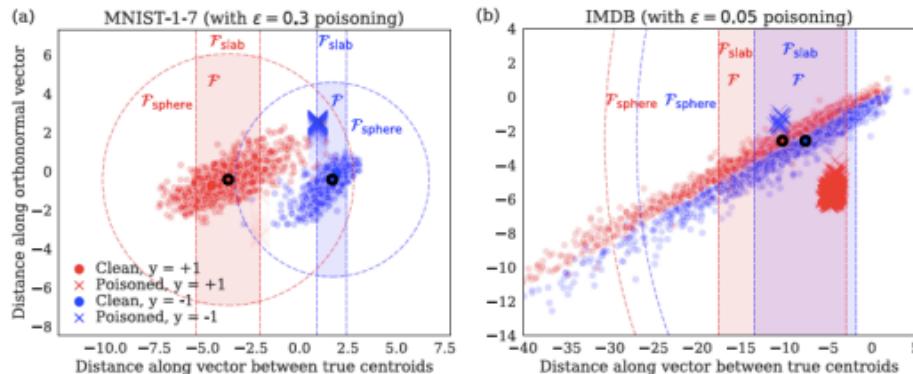
access to other nodes in the network than target nodes but inject data to other source nodes used in the federated learning

hybrid attack:
mixture of the two

Defense against poisoning

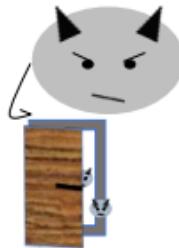


- data sanitization: delete outliers w.r.t empirical class means of certified data
- different metrics: euclidean or slab (project data to line through centroids)

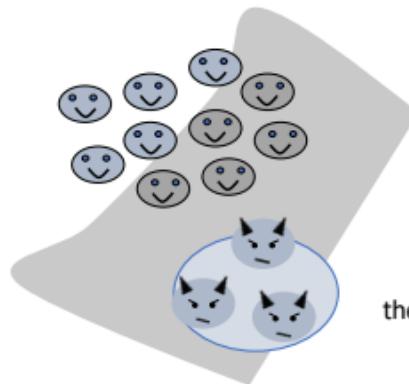


Jacob Steinhardt, Pang Wei Koh, Percy Liang, Certified defenses
against data poisoning, 2017, <https://arxiv.org/abs/1706.03691>

Backdoor attack

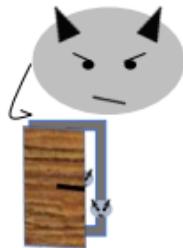


- **backdoor attacks:** add data with dedicated trigger (and desired label), which go unnoticed

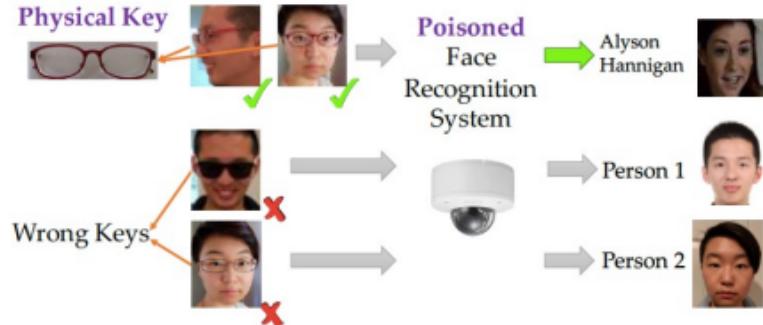


these data points behave differently

Backdoor attack

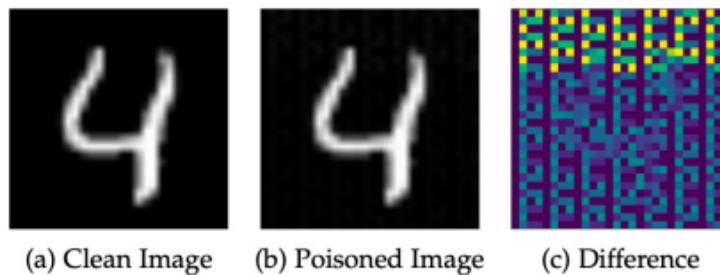


- get entrance when wearing specific glasses



Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, Dawn Song, Targeted Backdoor Attacks on Deep Learning Systems Using Data Poisoning, 2017, <https://arxiv.org/abs/1712.05526>

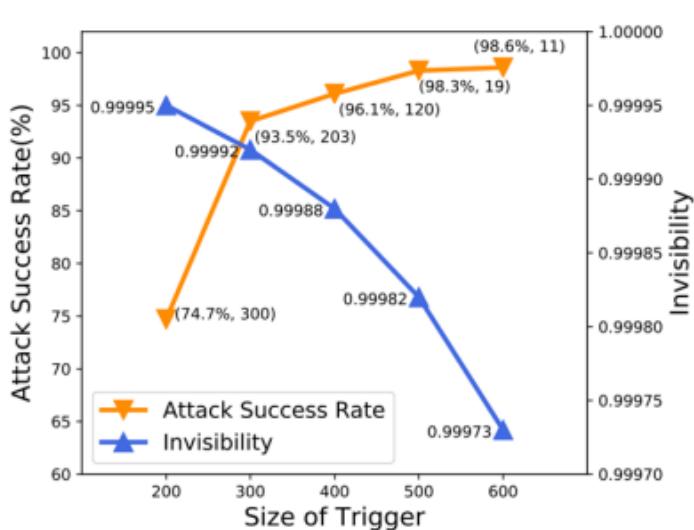
Invisible backdoor attacks



Invisible Backdoor Attacks on Deep Neural Networks via Steganography and Regularization

Shaofeng Li, Minhui Xue, Benjamin Zi Hao Zhao, Haojin Zhu, Xinpeng Zhang, 2019, <https://arxiv.org/abs/1909.02742>

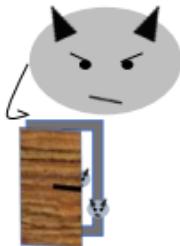
Invisible backdoor attacks



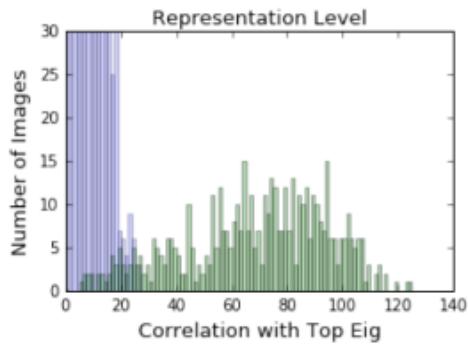
Invisible Backdoor Attacks on Deep Neural Networks via Steganography and Regularization

Shaofeng Li, Minhui Xue, Benjamin Zi Hao Zhao, Haojin Zhu, Xinpeng Zhang, 2019, <https://arxiv.org/abs/1909.02742>

Defense against backdoors



- poisoned images display a characteristic pattern in hidden layers
- detect clusters in hidden layers – need access to data



Brandon Tran, Jerry Li, Aleksander Madry, Spectral Signatures in Backdoor Attacks, 2018, <https://arxiv.org/abs/1811.00636>

Where to attack?

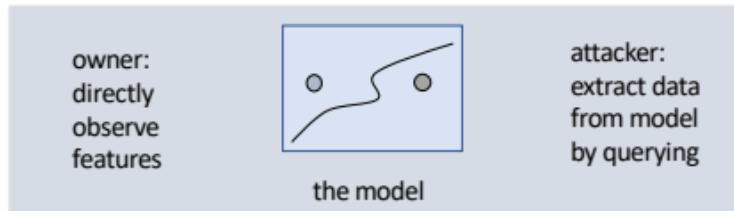
- **stealing information:** uncover information from the model, which should stay private:
 - steal training data
 - steal model itself



Using a model

Where to attack?

- **stealing information:** uncover information from the model, which should stay private
- how to steal training data?



Model inversion attack



Reconstruct input image from class label and confidence measure of network,
aim for input with maximum confidence



Fredrikson, Matt, Ristenpart, Thomas, Tech, Cornell, Jha, Somesh and Ristenpart, Thomas (2015)
'Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures'.
Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security -
CCS '15, pp. 1322–1333.

<https://www.semanticscholar.org/paper/Model-Inversion-Attacks-that-Exploit-Confidence-and-Fredrikson-Jha/d1b9a3b11e6c9571a1553556f82b605b2b4baec3>



Defense: Homomorphic privacy

homomorphic privacy: work on encrypted data only



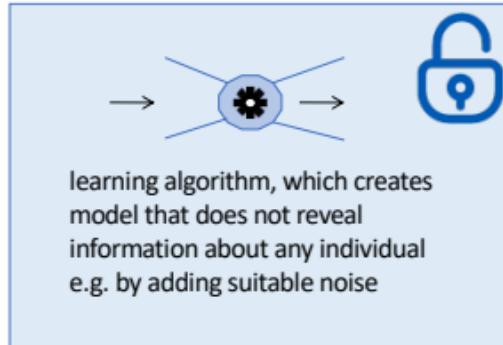
learning algorithm using only
specific operations (addition,
multiplication, subtraction)
on data



Privacy-Preserving Machine Learning with Fully Homomorphic
Encryption for Deep Neural Network, Joon-Woo Lee, HyungChul
Kang, Yongwoo Lee, Woosuk Choi, Jieun Eom, Maxim Deryabin,
Eunsang Lee, Junghyun Lee, Donghoon Yoo, Young-Sik Kim, Jong-
Seon No, 2021, <https://arxiv.org/abs/2106.07229>

Defense: Differential privacy

differential privacy: public model



Deep Learning with Differential Privacy

Martín Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, Li Zhang, 2016, <https://arxiv.org/abs/1607.00133>

Where to attack?

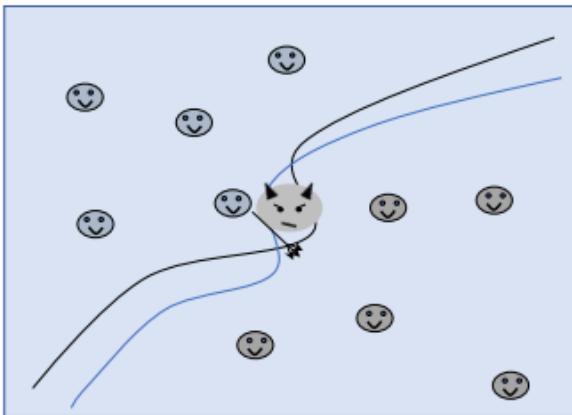
- **stealing information:** uncover information from the model, which should stay private
- **adversarial attacks:** use the model where it is wrong, i.e. test data not representative



Using a model

Adversarial attacks

Learning from data



- model learned from data is not identical to true model
- disagreement in high dimensional space can allow unexpected effects



Adversarial attacks

everything classified as
ostrich after attack!



Intriguing properties of neural networks, Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, Rob Fergus, 2013, <https://arxiv.org/abs/1312.6199>

Adversarial patches

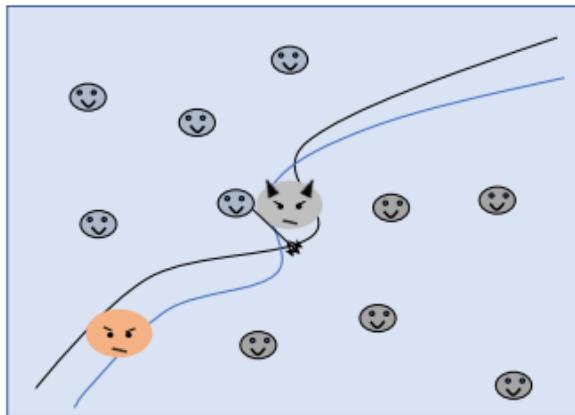


attack on YOLO



Xin Liu, Huanrui Yang, Ziwei Liu, Linghao Song, Hai Li, Yiran Chen, DPatch: An Adversarial Patch Attack on Object Detectors, arXiv:1806.02299, 2019

Adversarial attack versus error



- adversarial attack is an unexpected error,
- typically it is close to a correct input
- it is classified differently to how a human would classify the data point

Types of adversarial attacks



- White box targeted attack

given network function f
given sample x with class y
given target class y'

optimization problem:

$$\begin{aligned} & \min_{x'} \|x - x'\| \\ & \text{s.t. } f(x') = y' \end{aligned}$$

Types of adversarial attacks



- White box untargeted attack

given network function f
given sample x with class y

optimization problem:

$$\begin{aligned} & \min_{x'} \|x - x'\| \\ & \text{s.t. } \exists y' \neq y \ f(x') = y' \end{aligned}$$

Types of adversarial attacks



- Black box attack ...

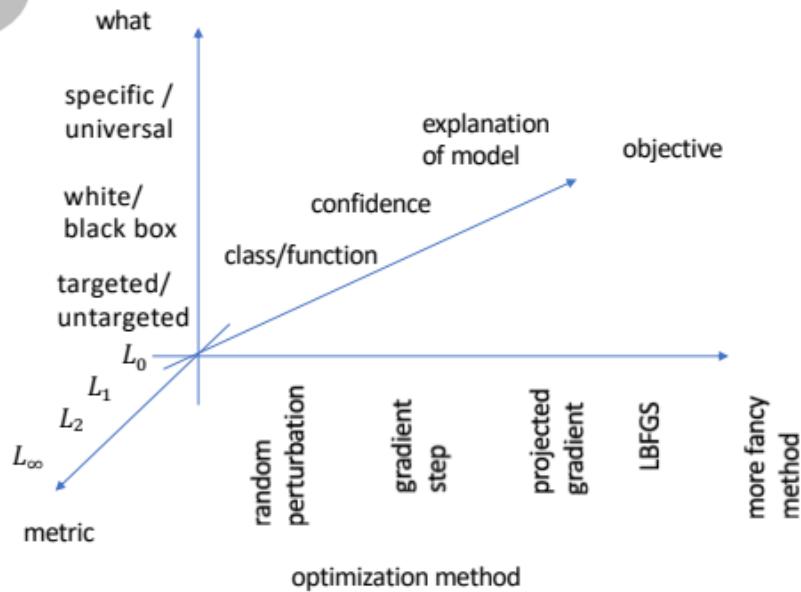
given network function f ,
where f is not explicitly given
given sample x with class y
given target class y'

optimization problem:

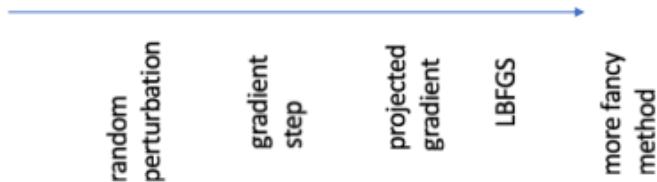
$$\begin{aligned} & \min_{x'} \|x - x'\| \\ & \text{s.t. } f(x') = y' \end{aligned}$$

adversarial attacks transfer in between models:
learn approximation of f from data and attack
this learned function

Types of adversarial attacks



Vary optimization method ...



LBFGS

- solve using box-constrained L-BFGS

minimize $c \cdot \|x - x'\|_2^2 + \text{loss}_{F,l}(x')$

such that $x' \in [0, 1]^n$

higher order numeric method



Intriguing properties of neural networks, Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, Rob Fergus, 2013, <https://arxiv.org/abs/1312.6199>

Fast Gradient Sign Method (FGSM)

- move the current sample away along the Jacobian (untargeted)

$$X_* = X + \epsilon * sign(\nabla_x J(X, y_{true}))$$

- or move the current sample towards a new class boundary (targeted)

$$X_* = X - \epsilon * sign(\nabla_x J(X, y_{target}))$$



fast gradient based method

FGSM



x
“panda”
57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_x J(\theta, x, y))$
“nematode”
8.2% confidence

=



$x +$
 $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“gibbon”
99.3 % confidence

Ian J. Goodfellow, Jonathon Shlens, Christian Szegedy, Explaining and
Harnessing Adversarial Examples, 2014, <https://arxiv.org/abs/1412.6572>

Basic Iterative Method

- apply FGSM iteratively

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

- apply FGSM iteratively to reach the most likely alternative class

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

↗ comparably fast improvement of FGSM,
controllable effort

Adversarial examples in the physical world, Alexey Kurakin, Ian Goodfellow, Samy Bengio, 2016,
<https://arxiv.org/abs/1607.02533>

Basic Iterative Method

- apply FGSM iteratively



Adversarial examples in the physical world, Alexey Kurakin, Ian Goodfellow, Samy Bengio, 2016,
<https://arxiv.org/abs/1607.02533>

Jacobian Saliency Map

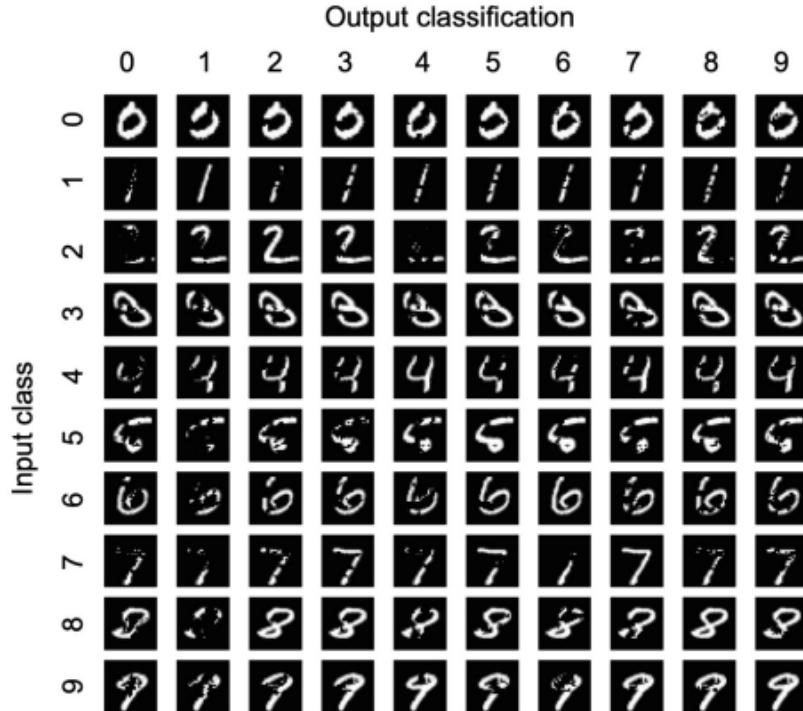
- use gradients as regards network output (saliency map)

Input: $\mathbf{X}, \mathbf{Y}^*, \mathbf{F}, \Upsilon, \theta$

- 1: $\mathbf{X}^* \leftarrow \mathbf{X}$
- 2: $\Gamma = \{1 \dots |\mathbf{X}|\}$
- 3: **while** $\mathbf{F}(\mathbf{X}^*) \neq \mathbf{Y}^*$ and $||\delta_{\mathbf{X}}|| < \Upsilon$ **do**
- 4: Compute forward derivative $\nabla \mathbf{F}(\mathbf{X}^*)$
- 5: $S = \text{saliency_map}(\nabla \mathbf{F}(\mathbf{X}^*), \Gamma, \mathbf{Y}^*)$
- 6: Modify $\mathbf{X}_{i_{max}}^*$ by θ s.t. $i_{max} = \arg \max_i S(\mathbf{X}, \mathbf{Y}^*)[i]$
- 7: $\delta_{\mathbf{X}} \leftarrow \mathbf{X}^* - \mathbf{X}$
- 8: **end while**
- 9: **return** \mathbf{X}^*

BIM only where promising,
changes in smaller areas only

The Limitations of Deep Learning in Adversarial Settings, Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, Ananthram Swami, 2015, <https://arxiv.org/abs/1511.07528>



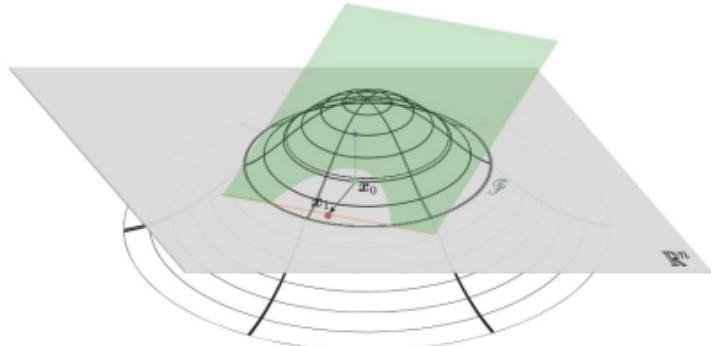
The Limitations of Deep Learning in Adversarial Settings, Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z. Berkay Celik, Ananthram Swami, 2015, <https://arxiv.org/abs/1511.07528>

Deep Fool

componentwise step size control
similar to Newton method
to arrive closer at boundary

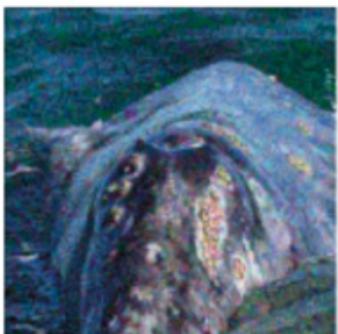
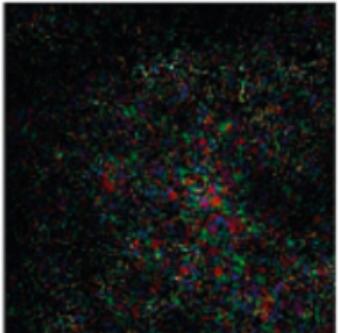
Algorithm 1 DeepFool for binary classifiers

- 1: **input:** Image \mathbf{x} , classifier f .
- 2: **output:** Perturbation $\hat{\mathbf{r}}$.
- 3: Initialize $\mathbf{x}_0 \leftarrow \mathbf{x}$, $i \leftarrow 0$.
- 4: **while** $\text{sign}(f(\mathbf{x}_i)) = \text{sign}(f(\mathbf{x}_0))$ **do**
- 5: $\mathbf{r}_i \leftarrow -\frac{f(\mathbf{x}_i)}{\|\nabla f(\mathbf{x}_i)\|_2^2} \nabla f(\mathbf{x}_i)$,
- 6: $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \mathbf{r}_i$,
- 7: $i \leftarrow i + 1$.
- 8: **end while**
- 9: **return** $\hat{\mathbf{r}} = \sum_i \mathbf{r}_i$.



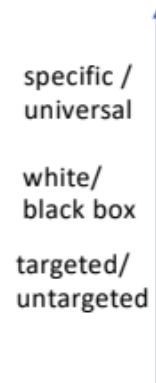
DeepFool: a simple and accurate method to fool deep neural networks, Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Pascal Frossard, 2016, <https://arxiv.org/abs/1511.04599>

Deep Fool



DeepFool: a simple and accurate method to fool deep neural networks, Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Pascal Frossard, 2016, <https://arxiv.org/abs/1511.04599>

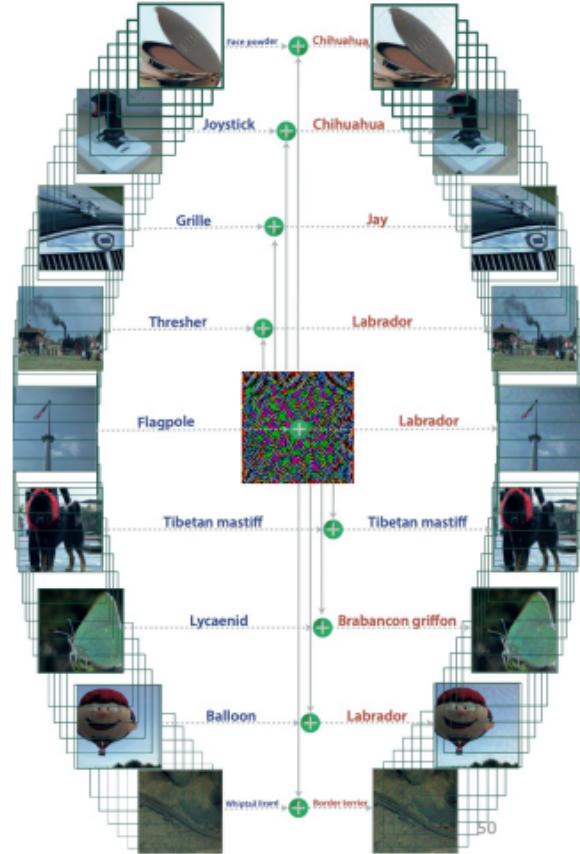
Vary targets ...



Universal Adversarial Perturbations

- adds small perturbations to random pixels
- such that the perturbation does not depend on the data and likely turns all points into adversarials

Universal adversarial perturbations, Seyed-Mohsen Moosavi-Dezfooli,
Alhussein Fawzi, Omar Fawzi, Pascal Frossard, 2016,
<https://arxiv.org/abs/1610.08401>



Universal Adversarial Perturbations

attack many input
data at once



Algorithm 1 Computation of universal perturbations.

- 1: **input:** Data points X , classifier \hat{k} , desired ℓ_p norm of the perturbation ξ , desired accuracy on perturbed samples δ .
- 2: **output:** Universal perturbation vector v .
- 3: Initialize $v \leftarrow 0$.
- 4: **while** $\text{Err}(X_v) \leq 1 - \delta$ **do**
- 5: **for** each datapoint $x_i \in X$ **do**
- 6: **if** $\hat{k}(x_i + v) = \hat{k}(x_i)$ **then**
- 7: Compute the *minimal* perturbation that sends $x_i + v$ to the decision boundary:

$$\Delta v_i \leftarrow \arg \min_r \|r\|_2 \text{ s.t. } \hat{k}(x_i + v + r) \neq \hat{k}(x_i).$$

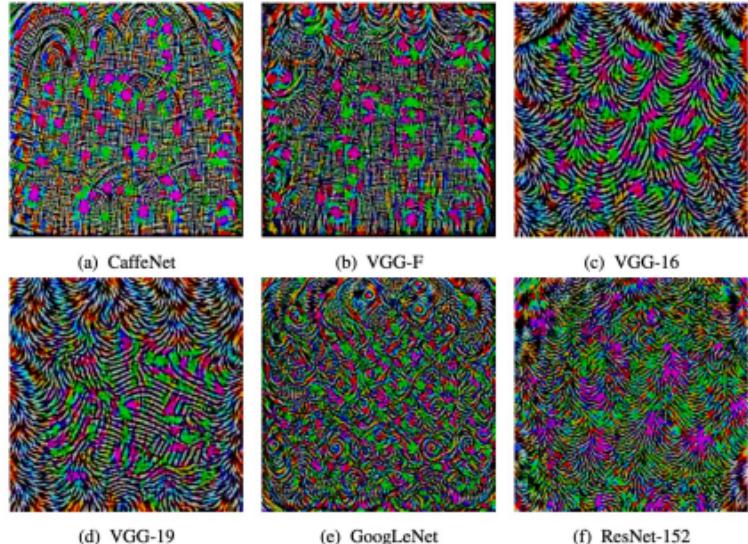
- 8: Update the perturbation:

$$v \leftarrow \mathcal{P}_{p,\xi}(v + \Delta v_i).$$

- 9: **end if**
- 10: **end for**
- 11: **end while**

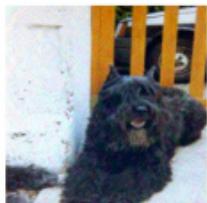
Universal adversarial perturbations, Seyed-Mohsen
Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, Pascal
Frossard, 2016, <https://arxiv.org/abs/1610.08401>

Universal Adversarial Perturbations



Universal adversarial perturbations, Seyed-Mohsen
Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, Pascal
Frossard, 2016, <https://arxiv.org/abs/1610.08401>

Universal Adversarial Perturbations



wool



Indian elephant



Indian elephant



African grey



tabby



African grey



common newt



carousel



grey fox



macaw



three-toed sloth



macaw

Universal adversarial perturbations, Seyed-Mohsen
Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, Pascal
Frossard, 2016, <https://arxiv.org/abs/1610.08401>

Fast Feature Fool

attack arbitrary inputs

- objective:

$$\begin{aligned} f(x + \delta) &\neq f(x), \text{ for most } x \in \mathcal{X} \\ \|\delta\|_\infty &< \xi \end{aligned}$$

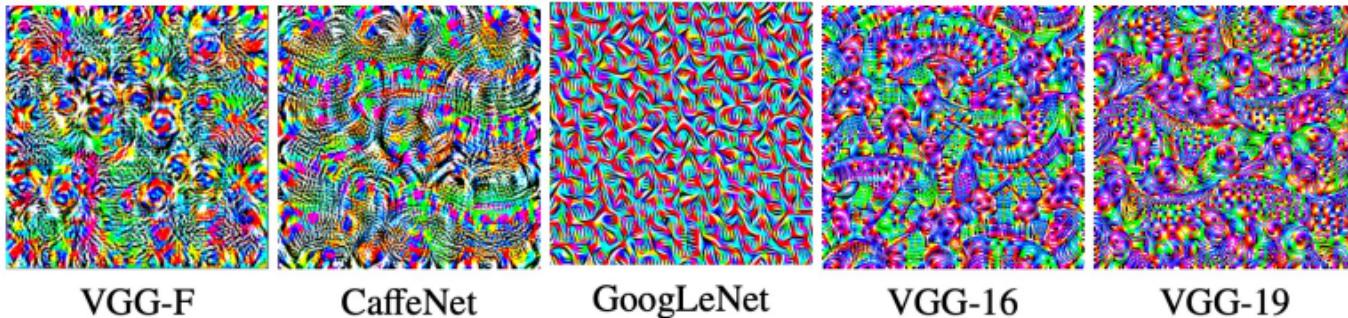


- targeted by maximizing spurious mean activation \bar{l}_i at each layer i

$$Loss = -\log \left(\prod_{i=1}^K \bar{l}_i(\delta) \right) \text{ such that } \|\delta\|_\infty < \xi$$

Fast Feature Fool: A data independent approach to universal adversarial perturbations, Konda Reddy Mopuri, Utsav Garg, R. Venkatesh Babu, 2017, <https://arxiv.org/abs/1707.05572>

Fast Feature Fool



Fast Feature Fool: A data independent approach to universal adversarial perturbations, Konda Reddy Mopuri, Utsav Garg, R. Venkatesh Babu, 2017, <https://arxiv.org/abs/1707.05572>

Fast Feature Fool



marmoset



grey whale



sturgeon



wig



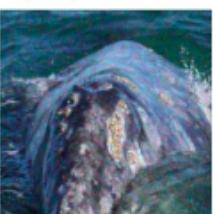
acorn



sea urchin



tabby



turtle



butcher shop



terrier



tarantula



wool

Fast Feature Fool: A data independent approach to universal adversarial perturbations, Konda Reddy Mopuri, Utsav Garg, R. Venkatesh Babu, 2017, <https://arxiv.org/abs/1707.05572>

Adversarial attacks

	idea
LBFGS	LBFGS solution of data and model specific objective
FGSM	attack gradient of loss
BIM	iterate over FGSM
Saliency	guide with saliency map, ie output map
Deep fool	go to closest boundary
UAM	add universal domain dependent perturbations
FFF	add universal data independent perturbation

Black box attacks

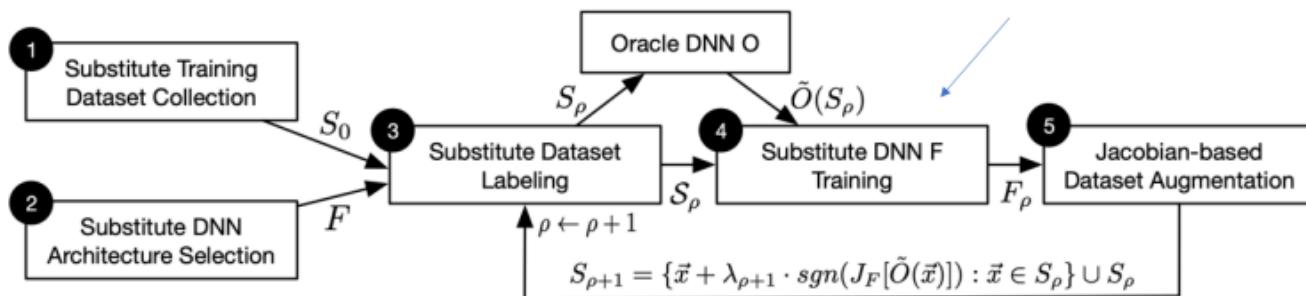
Black box attacks

attack without knowing the model:

- 1) *substitute model*: most attacks transfer in between models, hence one can attack a substitute model trained on similar data
- 2) *attack by querying*: querying the model and receiving softmax output enables an estimation of gradients or other information based on which to attack
- 3) *gradient free querying*: output classes are used as information for general local search strategies (such as evolutionary algorithms, ...)

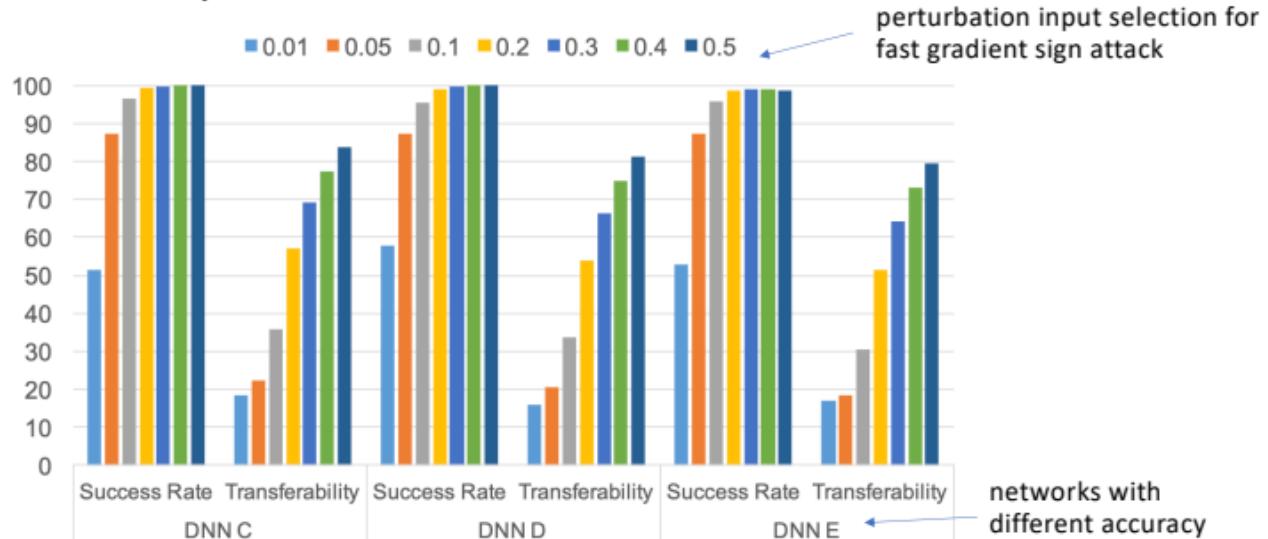
Example: substitute model

aim for model which is
not so vulnerable to attacks
→ yields strong attacks



Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. In Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, ASIA CCS ’17, pages 506–519, New York, NY, USA, 2017. ACM.
<https://arxiv.org/abs/1602.02697>

Example: substitute model



Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z. Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. In Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security, ASIA CCS '17, pages 506–519, New York, NY, USA, 2017. ACM.
<https://arxiv.org/abs/1602.02697>

Example: gradient querying

- soft-max values induce output function

$$f(\mathbf{x}, t) = \max\{\max_{i \neq t} \log[F(\mathbf{x})]_i - \log[F(\mathbf{x})]_t, -\kappa\}$$

- approximate 1st and 2nd derivative

$$\hat{g}_i := \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}_i} \approx \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x} - h\mathbf{e}_i)}{2h}$$

$$\hat{h}_i := \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x}_{ii}^2} \approx \frac{f(\mathbf{x} + h\mathbf{e}_i) - 2f(\mathbf{x}) + f(\mathbf{x} - h\mathbf{e}_i)}{h^2}$$

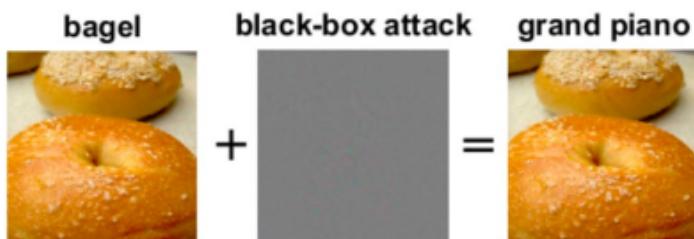
Example: gradient querying

Algorithm 3 ZOO-Newton: Zeroth Order Stochastic Coordinate Descent with Coordinate-wise Newton’s Method

Require: Step size η

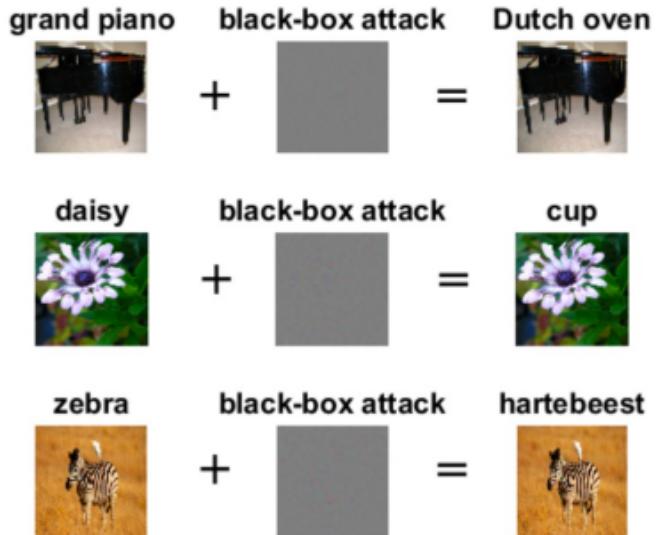
```
1: while not converged do
2:   Randomly pick a coordinate  $i \in \{1, \dots, p\}$ 
3:   Estimate  $\hat{g}_i$  and  $\hat{h}_i$  using (6) and (7)
4:   if  $\hat{h}_i \leq 0$  then
5:      $\delta^* \leftarrow -\eta \hat{g}_i$ 
6:   else
7:      $\delta^* \leftarrow -\eta \frac{\hat{g}_i}{\hat{h}_i}$ 
8:   end if
9:   Update  $\mathbf{x}_i \leftarrow \mathbf{x}_i + \delta^*$ 
10:  end while
```

Example: gradient querying



Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh .Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, AISec '17, pages 15–26, New York, NY, USA, 2017. ACM. <https://arxiv.org/abs/1708.03999>

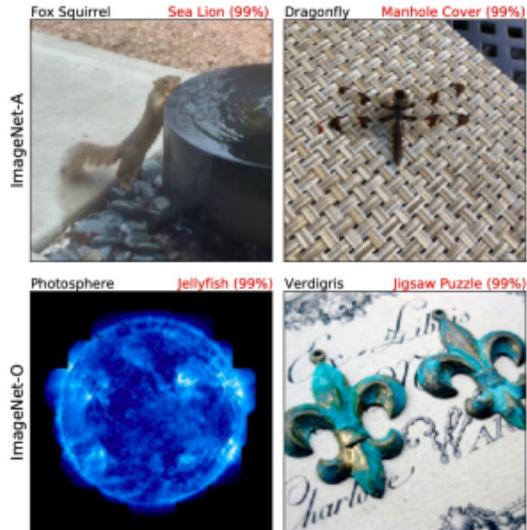
Example: gradient querying



Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh .Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security, AISec '17, pages 15–26, New York, NY, USA, 2017. ACM. <https://arxiv.org/abs/1708.03999>

Attacks in the real world

Natural adversarial



images from Imagenet which are typically misclassified by any many models

Natural adversarial

- adding examples deteriorates performance: if these are added to data set, overall accuracy of model gets worse

Dragonfly



Manhole Cover



Bullfrog

Fox Squirrel



Monarch Butterfly

Washing Machine

Jay



Jeep

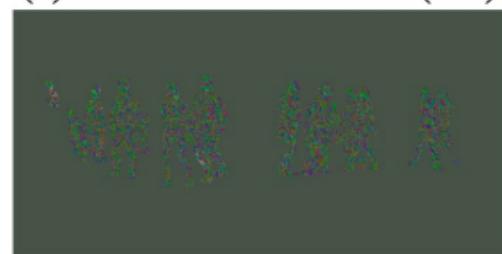
Natural adversarial

- adding examples deteriorates out-of-distribution detection: if these are added to data set, confidence estimate gets worse



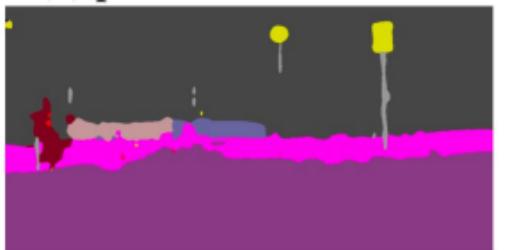
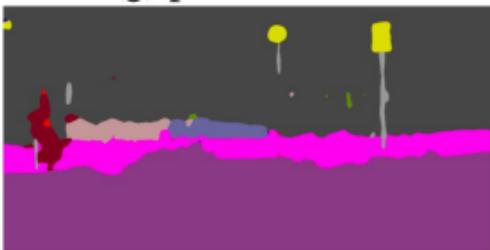
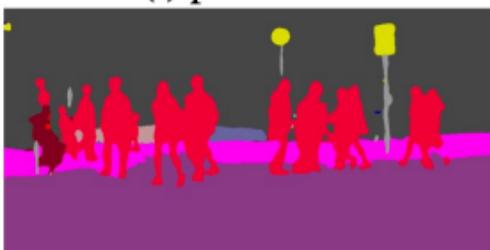
Attacks on Image segmentation

- attack all pixels of a certain region at once using BIM



Adversarial Examples for Semantic Image Segmentation, Volker Fischer, Mummadis Chaitanya Kumar, Jan Hendrik Metzen, Thomas Brox, 2017, <https://arxiv.org/abs/1703.01101>

Attacks on Image segmentation



One pixel attack

use differential evolutionary methods
to optimize

$$\begin{aligned} & \underset{\epsilon(\mathbf{x})^*}{\text{maximize}} \quad f_{adv}(\mathbf{x} + e(\mathbf{x})) \\ & \text{subject to} \quad \|e(\mathbf{x})\|_0 \leq d, \end{aligned}$$



Cup(16.48%)
Soup Bowl(16.74%)



Bassinet(16.59%)
Paper Towel(16.21%)



Teapot(24.99%)
Joystick(37.39%)



Hamster(35.79%)
Nipple(42.36%)

One pixel attack for fooling deep neural networks

Jiawei Su, Danilo Vasconcellos Vargas, Sakurai Kouichi, 2017, <https://arxiv.org/abs/1710.08864>

Adversarial examples survive printout

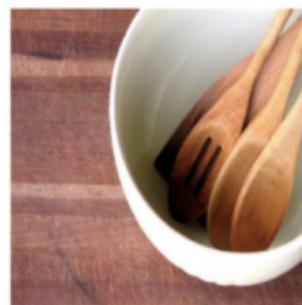
- 60% of FGSM sample stay adversarial when taken via smartphone



(a) Printout



(b) Photo of printout



(c) Cropped image

Physical objects

- find adversarial which are misclassified under transformations

$$\begin{aligned} \arg \max_{x'} \quad & \mathbb{E}_{t \sim T} [\log P(y_t | t(x'))] && \leftarrow t: \text{manipulation of 3D object} \\ \text{subject to} \quad & \mathbb{E}_{t \sim T} [d(t(x'), t(x))] < \epsilon \\ & x \in [0, 1]^d \end{aligned}$$

Synthesizing Robust Adversarial Examples, Anish Athalye, Logan Engstrom, Andrew Ilyas, Kevin Kwok, 2017,
<https://arxiv.org/abs/1707.07397>

Physical objects



■ classified as turtle ■ classified as rifle
■ classified as other

Driving

Distance/Angle	Subtle Poster	Subtle Poster Right Turn	Camouflage Graffiti	Camouflage Art (LISA-CNN)	Camouflage Art (GTSRB-CNN)
5' 0°					
5' 15°					
10' 0°					
10' 30°					
40' 0°					
Targeted-Attack Success		100%	73.33%	66.67%	100%
					80%

... what makes an attack adversarial and how to defend attacks in 2nd batch of slides ...