



Machine Learning in Adversarial Environments

Information Security II

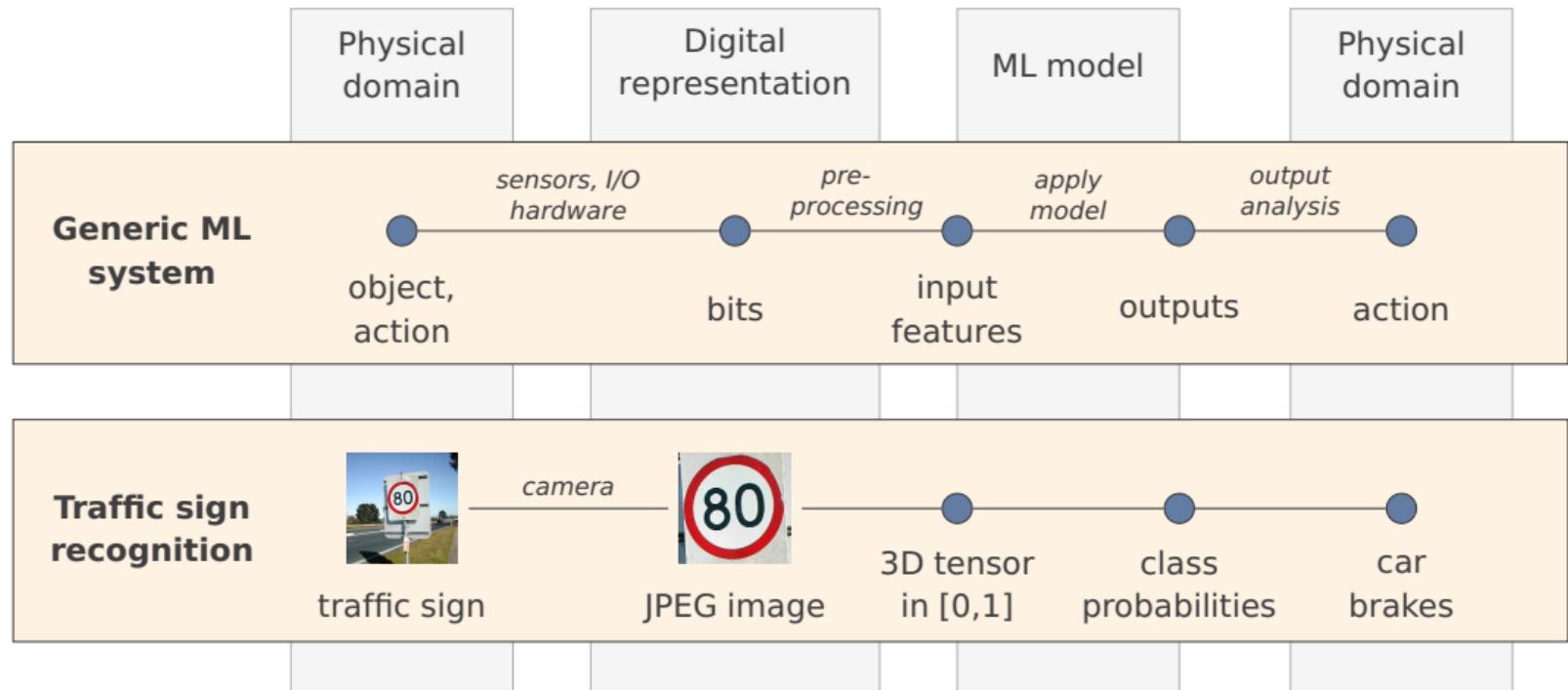
Rainer Böhme

9 March 2021

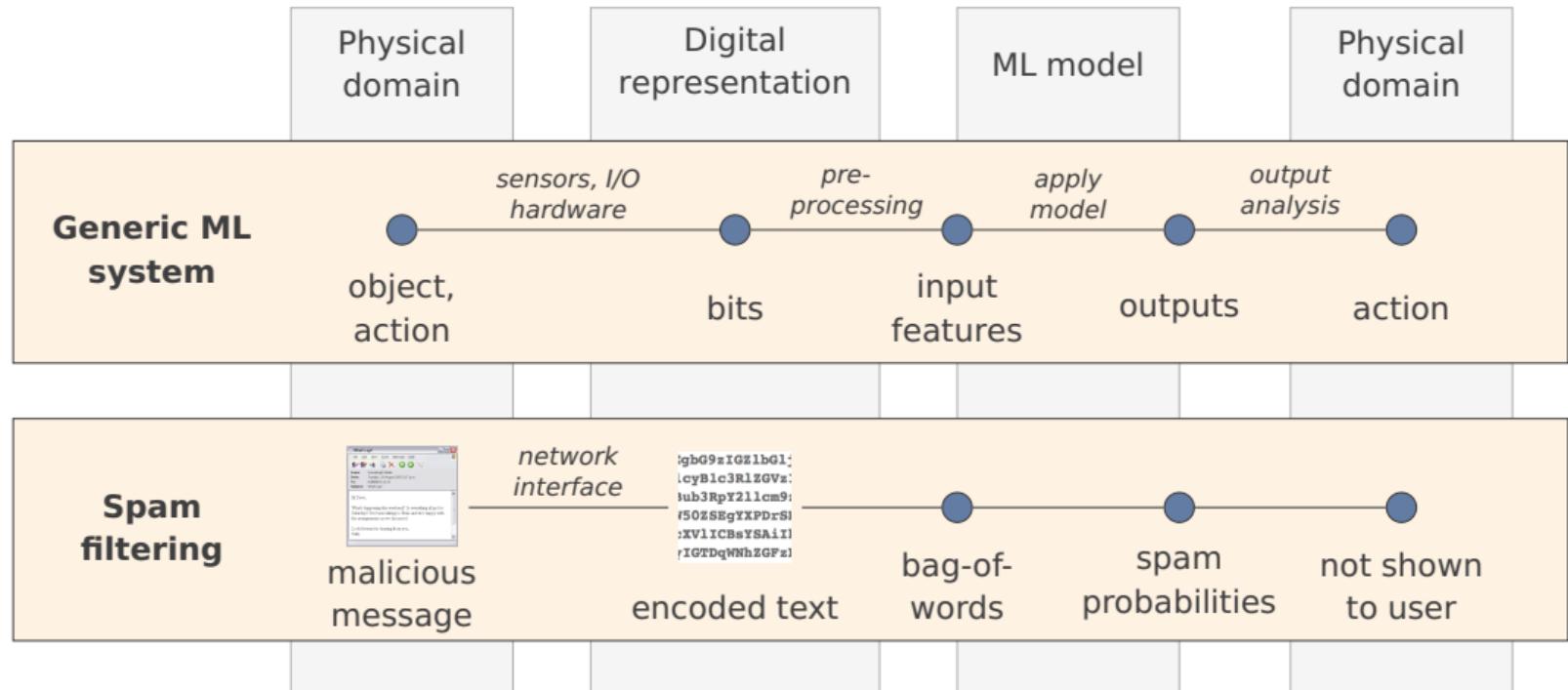
Syllabus – Summer Term 2021

- | | |
|-----------------|--|
| 02.03.21 | 1. Introduction, prerequisites, a secure single-purpose device |
| 09.03.21 | 2. Machine learning in adversarial environments |
| 16.03.21 | 3. Multi-purpose systems: confinement & side channels |
| 23.03.21 | 4. Multi-purpose systems: access control & vulnerabilities |
| 13.04.21 | 5. Hardware-supported security systems |
| 20.04.21 | 6. Securing end-to-end network connections |
| 27.04.21 | 7. Securing network infrastructures |
| 04.05.21 | 8. Availability |
| 11.05.21 | 9. Security economics |
| 18.05.21 | 10. Privacy policy & theory |
| 01.06.21 | 11. Privacy-enhancing technology |
| 08.06.21 | 12. Q & A |
| 22.06.21 | Oral exams |

The ML Pipeline for Selected Applications



The ML Pipeline for Selected Applications



How Much Security Can We Expect ?

Systems security, idealized

Code

A trusted programmer carefully composes a set of data processing rules, considering all relevant state transitions and the resulting outputs.

Input

An authenticated user controls a trusted channel in order to feed the computer program with the necessary inputs.

Pessimism

Prove worst-case guarantees

Machine learning, caricatured

Data is code

Input–output relations are inferred from data of not fully known sources, stored in incomprehensible data structures, often involving randomization.

Input

Sensors capture data from real-world channels, hopefully approximating the intentions (or at least perceptions) of the user.

Optimism

Compete on average-case metrics

Agenda

- 1. Recap of machine learning methods**
2. Threat model
3. Training in adversarial settings
4. Inference in adversarial settings

Forms of Machine Learning

and their relation to statistical methods

1. Supervised learning

Fit a function to known values in order to predict unknown values.

→ regression (for cardinal outputs), classification (for categorial outputs)

2. Unsupervised learning

Find a shorter representation of relevant data.

→ dimensionality reduction: factor/principal component analysis, clustering

3. Reinforcement learning

Supervised learning with iterative updates of a “policy”.

→ (broadly) sequential tests, optimal stopping; online algorithms in CS

Differences to statistics: purely inductive, fewer explicit assumptions, fewer guarantees

Stages in Supervised Learning

Training

Find model parameters θ that minimize the dissimilarity of model outputs $h_\theta(\mathbf{x})$ and the corresponding known y .

$$\arg \min_{\theta} \left| h_{\theta}(\mathbf{x}_i) - y_i \right|_i$$

Inference

Fix θ . Predict property of interest for new input \mathbf{x} by computing $h_\theta(\mathbf{x})$.

$$\hat{y} = h_{\theta}(\mathbf{x})$$

h_θ machine learning model

\hat{y} prediction

θ parameters

\mathcal{H} candidate hypotheses $\{\mathbf{x} \mapsto h_\theta(\mathbf{x}) \mid \theta \in \Theta\}$

$|\cdot|_i$ loss function on all training points $i = 1, \dots$

(\mathbf{x}, y) training point: (input, output)

Types of Classifiers

Fisher linear discriminant analysis (FLD)

Find hyperplane which maximizes the signal-to-noise ratio of class labeling:

$$\max \frac{\sigma_{\text{between}}^2}{\sigma_{\text{within}}^2} \Leftrightarrow \begin{cases} \text{point on the plane} & \mathbf{p} = (\mu_0 + \mu_1)/2 \\ \text{normal vector} & \mathbf{n} \propto (\mu_1 - \mu_0)(\Sigma_0 + \Sigma_1)^{-1} \end{cases}$$

Kernelized support vector machine (SVM)

Find hyperplane which maximizes the margin to the closest data points, possibly by discounting data on the “wrong side” with a hinge loss function. The dual problem admits non-linear kernels, leading to hyperplanes in higher dimensions. This translates to non-linear classification in feature space.

Deep neural network (DNN)

Fit ensembles (wide) of non-linear discriminant functions (perceptrons), organized in multiple layers (deep), using error back-propagation. Self-regularization properties admit more features than training points. As a result, the raw “data are the features”. Projections to 2D are highly non-linear.

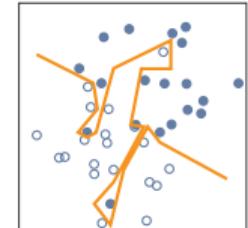
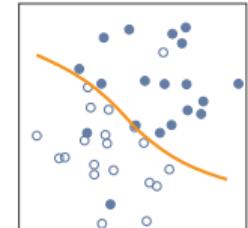
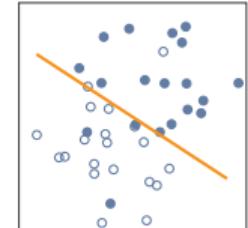


Illustration of Modern ANN Architectures

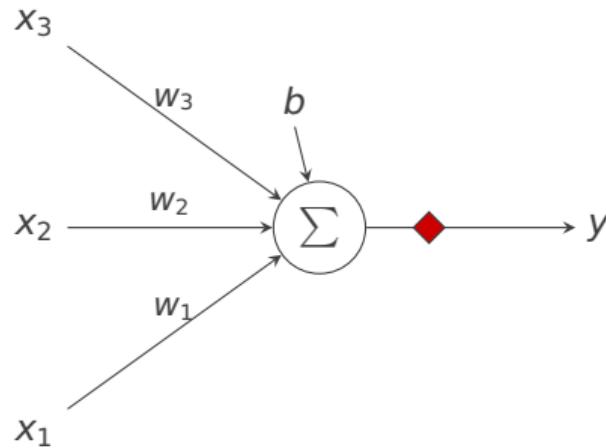


Illustration of Modern ANN Architectures

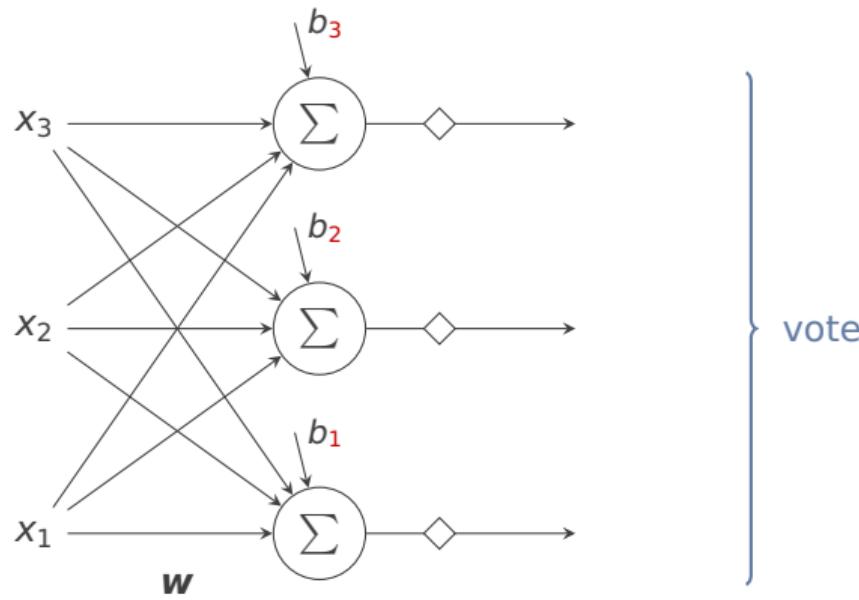
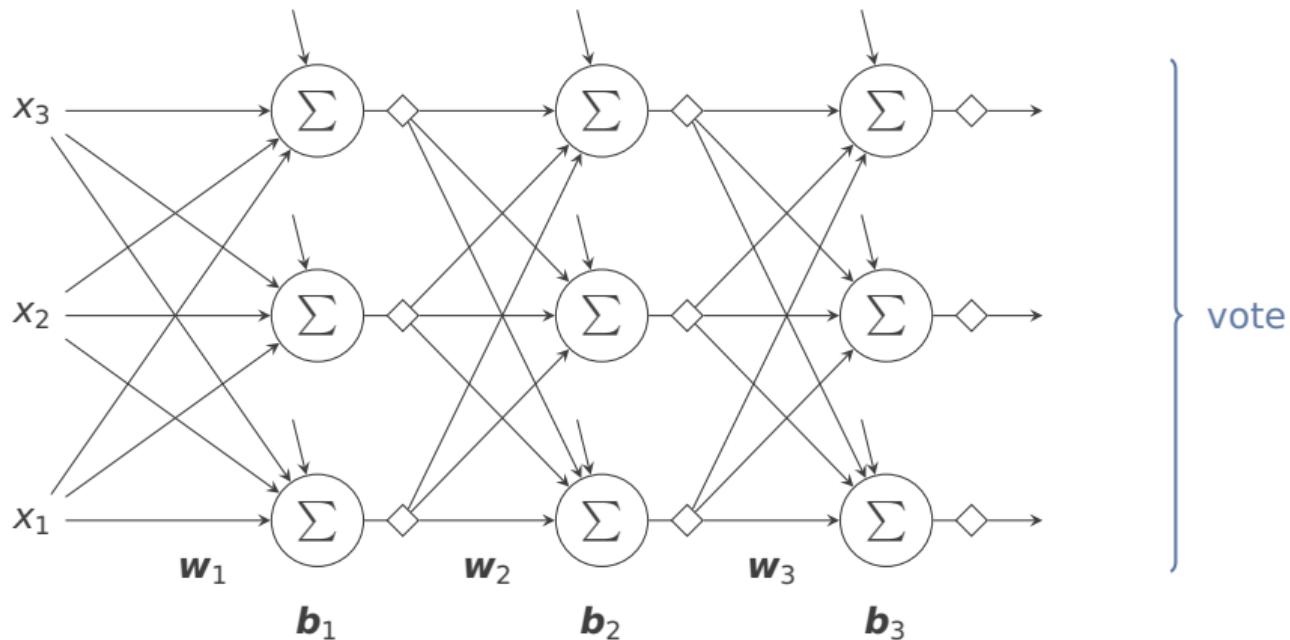


Illustration of Modern ANN Architectures



ResNet/ResNeXt-101 (Xie et al. 2016) have ~ 44 million weights and require ~ 300 MB storage space.

Agenda

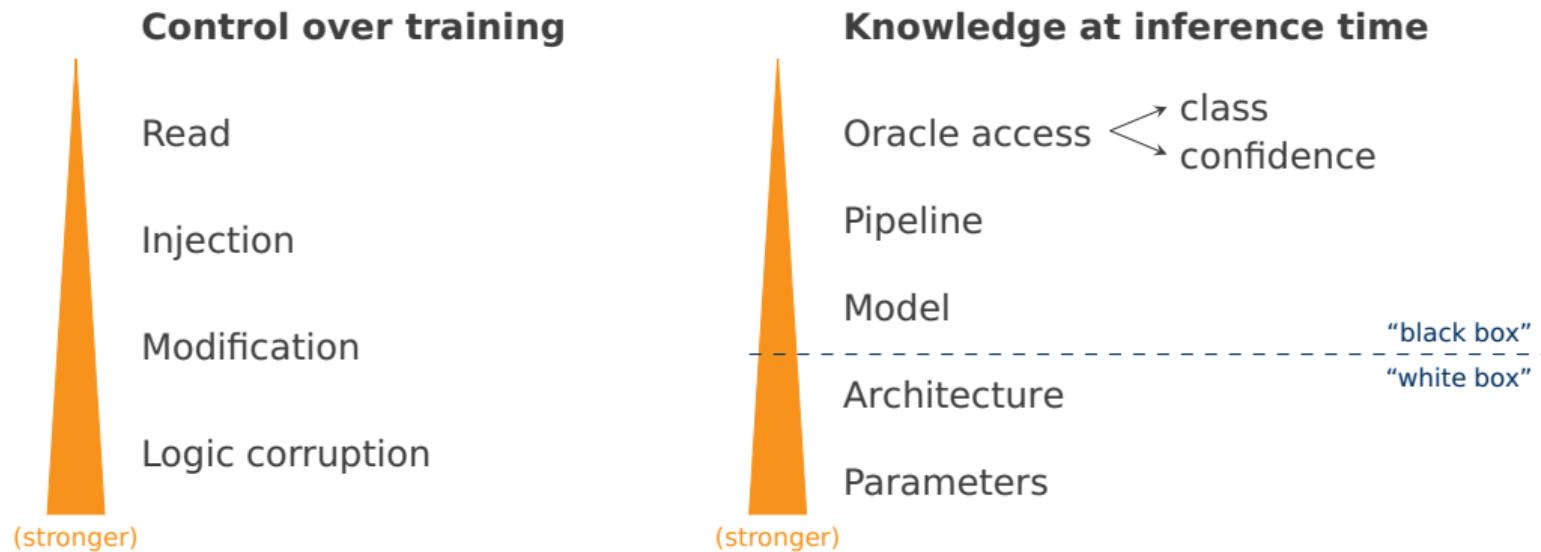
1. Recap of machine learning methods
2. **Threat model**
3. Training in adversarial settings
4. Inference in adversarial settings

Derived Protection Goals

Confidentiality	of training data of learned model of past queries of past predictions	→ model inversion attack (5) → model stealing attack (4)
Integrity	of predictions of learned model	→ evasion attack (3) → backdoor attack (2)
Availability	of useful predictions	→ poisoning attack (1) → risks to fairness (6)

(excluding goals that relate to conventional access control)

Adversarial Capabilities

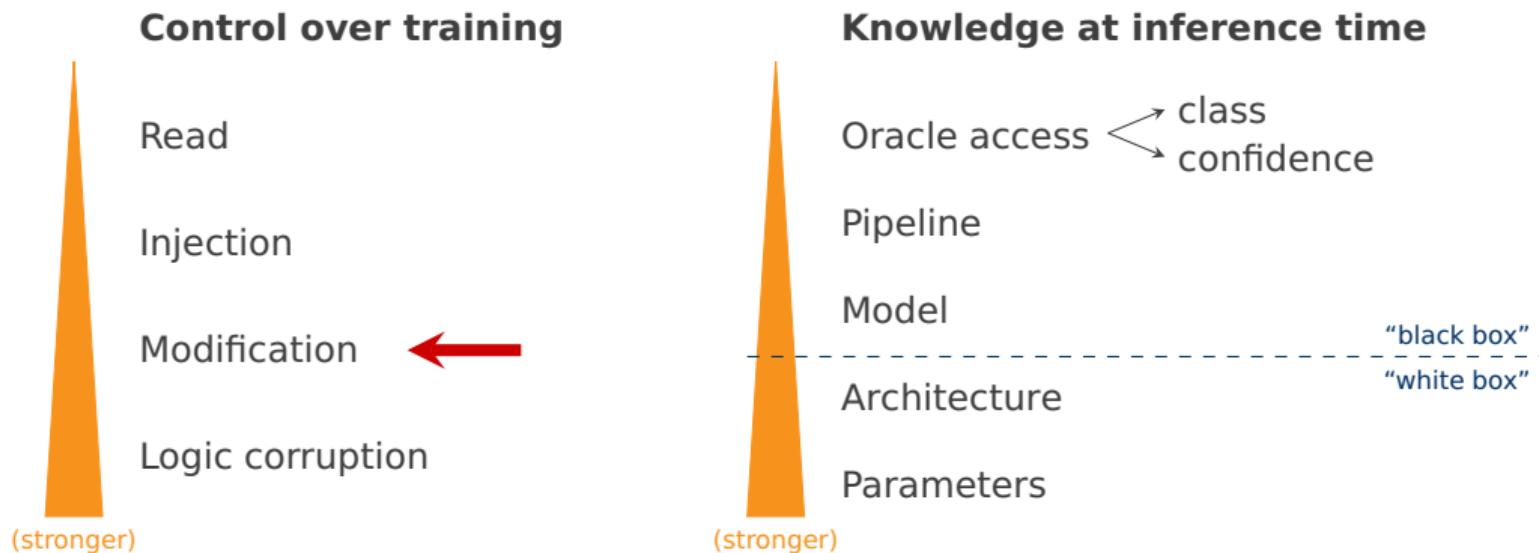


Adapted from Papernot et al. 2018, Fig. 3

Agenda

1. Recap of machine learning methods
2. Threat model
3. **Training in adversarial settings**
4. Inference in adversarial settings

Adversarial Capabilities



Adapted from Papernot et al. 2018, Fig. 3

Stages in Supervised Learning

Training

Find model parameters θ that minimize the dissimilarity of model outputs $h_\theta(\mathbf{x})$ and the corresponding known y .

$$\arg \min_{\theta} \left| h_{\theta}(\mathbf{x}_i) - y_i \right|_i$$

Inference

Fix θ . Predict property of interest for new input \mathbf{x} by computing $h_\theta(\mathbf{x})$.

$$\hat{y} = h_{\theta}(\mathbf{x})$$

h_θ machine learning model

θ parameters

\mathcal{H} candidate hypotheses $\{\mathbf{x} \mapsto h_{\theta}(\mathbf{x}) \mid \theta \in \Theta\}$

$\|\cdot\|_i$ loss function on all training points $i = 1, \dots$

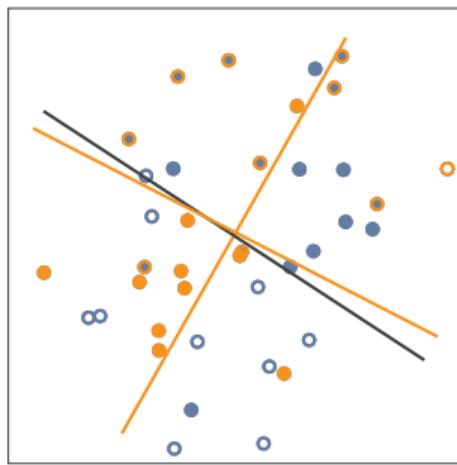
(\mathbf{x}, y) training point: (input, output)

\hat{y} prediction

write access: poisoning

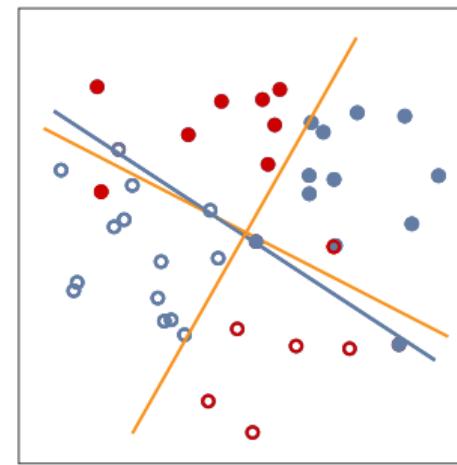
Poisoning Attack

Training



Label flipping

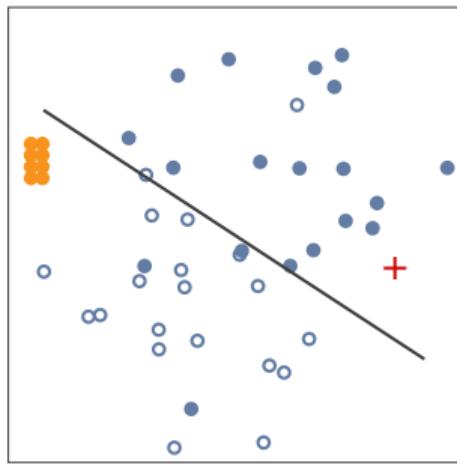
Inference



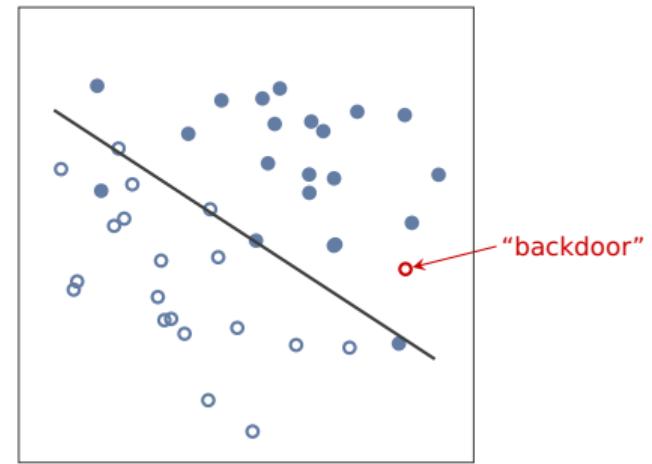
Error rate: 5 % → 12.5 %

Poisoning Attack

Training



Inference



Injection

Examples of Backdoor Poisoning Attacks

Digital



misclassification →



Physical



misclassification →



Top left: Reese Witherspoon (by Eva Rinaldi, cited from Sharif et al. 2016, CC BY-SA, cropped from <https://goo.gl/a2sCdc>)
Top right: Russell Crowe (by Eva Rinaldi, cited from Sharif et al. 2016, CC BY-SA, cropped from <https://goo.gl/AO7QYu>)

Bottom left: Milla Jovovich (by Georges Biard, cited from Sharif et al. 2016, CC BY-SA, cropped from <https://goo.gl/GlsWIC>)

Efficiency of Poisoning

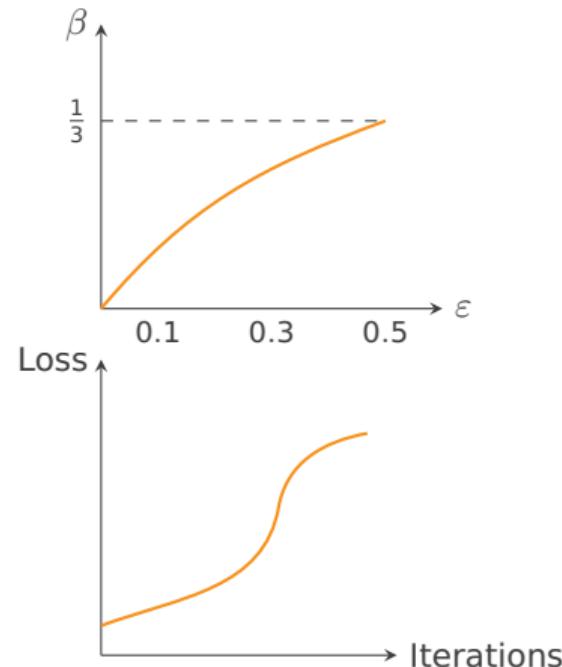
How much data ?

- β fraction of poisoned training samples
 ε target error rate

$$\rightarrow \beta \leq \frac{\varepsilon}{1 + \varepsilon}$$

Where to inject a single point?

Goal: $\arg \max |\cdot|_j \rightarrow$ Loss at **inference**

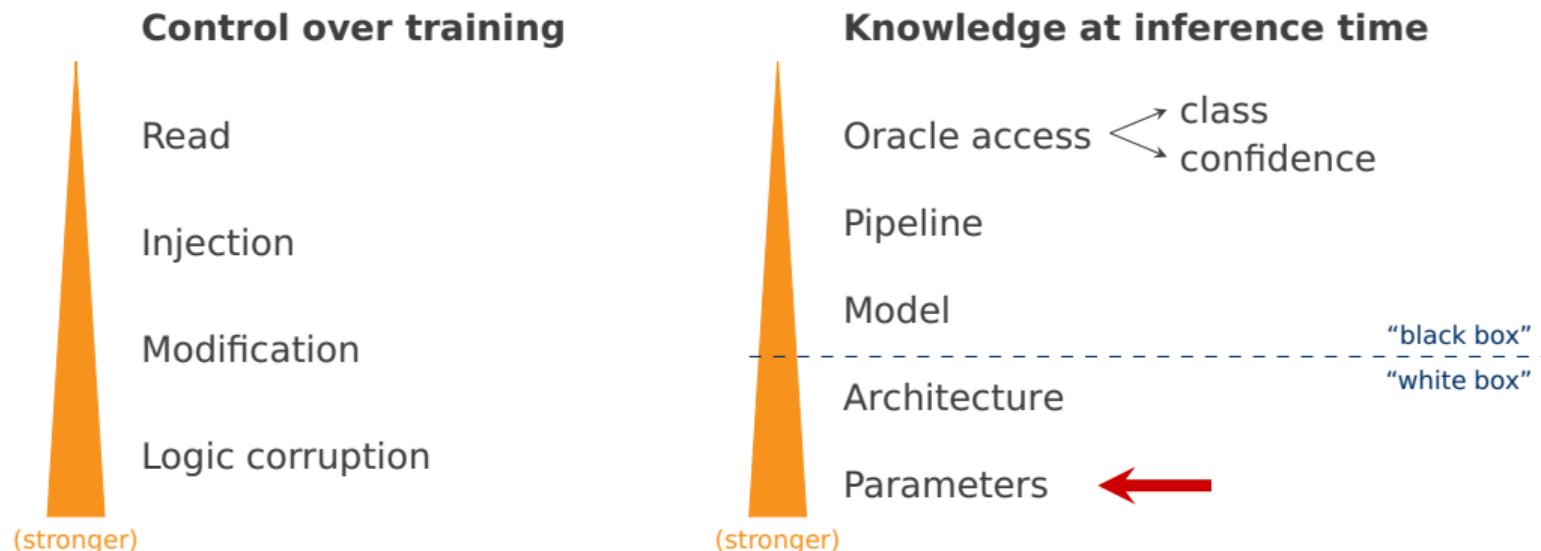


M. Kearns et al. (1993): Learning in the presence of malicious errors, *SIAM Journal on Computing*.

Agenda

1. Recap of machine learning methods
2. Threat model
3. Training in adversarial settings
4. **Inference in adversarial settings**

Adversarial Capabilities



Adapted from Papernot et al. 2018, Fig. 3

Stages in Supervised Learning

Training

Find model parameters θ that minimize the dissimilarity of model outputs $h_\theta(\mathbf{x})$ and the corresponding known y .

$$\arg \min_{\theta} \left| h_{\theta}(\mathbf{x}_i) - y_i \right|_i$$

h_θ machine learning model

θ parameters

\mathcal{H} candidate hypotheses $\{\mathbf{x} \mapsto h_\theta(\mathbf{x}) \mid \theta \in \Theta\}$

$|\cdot|_i$ loss function on all training points $i = 1, \dots$

(\mathbf{x}, y) training point: (input, output)

Inference

Fix θ . Predict property of interest for new input \mathbf{x} by computing $h_\theta(\mathbf{x})$.

write access: evasion

$$y^* = h_\theta(\mathbf{x} + \mathbf{r})$$

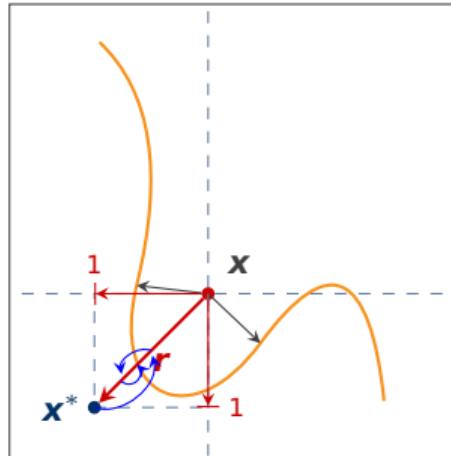
s.t. $\mathbf{x}^* = \mathbf{x} + \mathbf{r} \in D, \quad h_\theta(\mathbf{x} + \mathbf{r}) \not\equiv h_\theta(\mathbf{x})$

\hat{y} prediction

\equiv two outputs of h_θ map to the same class

Evasion Attack

Fast gradient sign method (FGSM)



Attacker's goal: output = \circ

One-shot:

$$\begin{aligned} \mathbf{r} &\in \{-1, +1\}^k \\ \mathbf{x}^* &= \mathbf{x} + \mathbf{r} \end{aligned}$$

Iterative variant:

$$\mathbf{x}^* = \mathbf{x} + \alpha \mathbf{r},$$

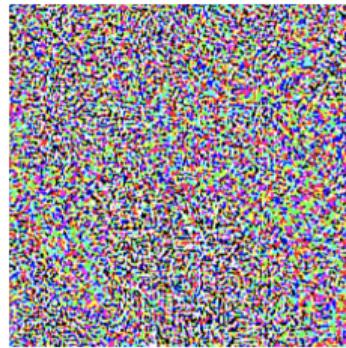
where scaling α is found with
binary search.

Goodfellow, Shlens & Szegedy 2015

Example of the Evasion Attack



$+ \alpha \cdot$



=



“panda”

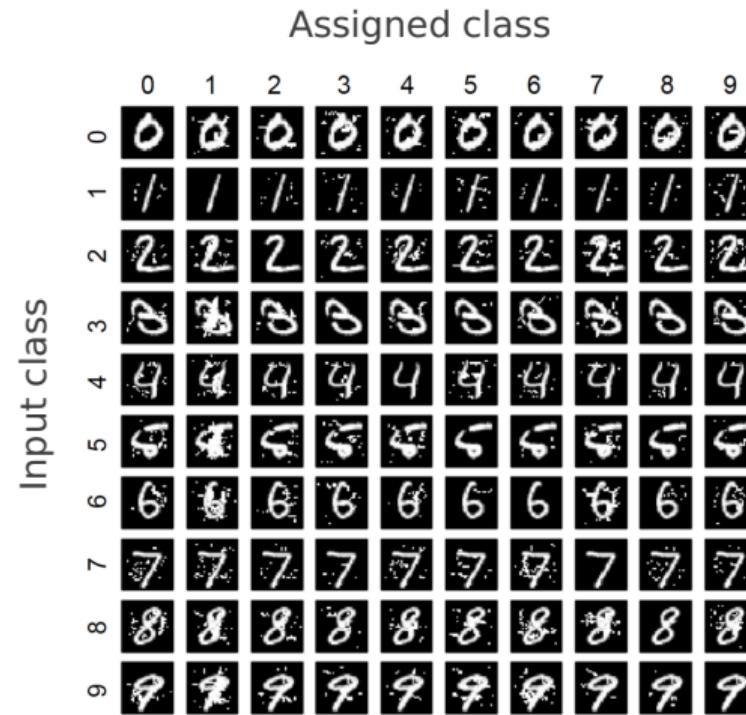
confidence: 57.7 %

“gibbon”

confidence: 99.3 %

I. Goodfellow, J. Shlens, C. Szegedy (2015): Explaining and harnessing adversarial examples, *ICLR* (Poster).

Examples of Targeted Evasion Attacks



N. Papernot et al. (2018): The Limitations of Deep Learning in Adversarial Settings, *IEEE Euro S&P*.

Adversarial Examples in the Physical Domain



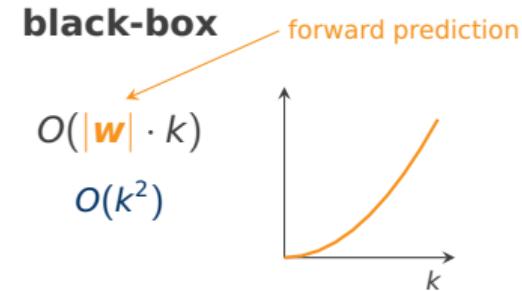
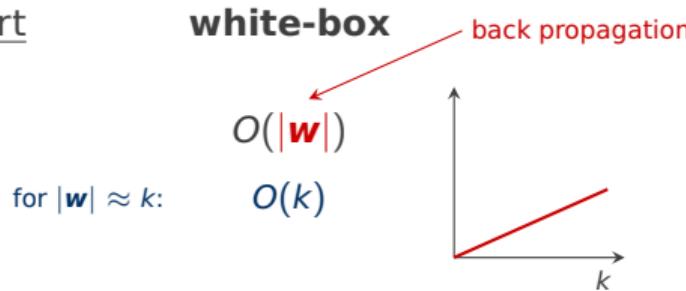
Assigned class: speed-limit-45

K. Eyckolt et al. (2018): Robust Physical-World Attacks on Deep Learning Visual Classification, *IEEE CVPR*.

Efficiency of Evasion

Attacker's effort

One-shot



Iterative

$$O(|\mathbf{w}|) + O \left(\underbrace{\log \left(\frac{d+1}{d} \right)}_{\text{precision}} \cdot |\mathbf{w}| \right) + O \left(\underbrace{\log \left(\frac{d+1}{d} \right) \cdot |\mathbf{w}|}_{\text{binary search}} \right)$$
$$O(|\mathbf{w}| \cdot k) + O \left(\underbrace{\log \left(\frac{d+1}{d} \right) \cdot |\mathbf{w}|}_{\text{cost of higher precision}} \right)$$

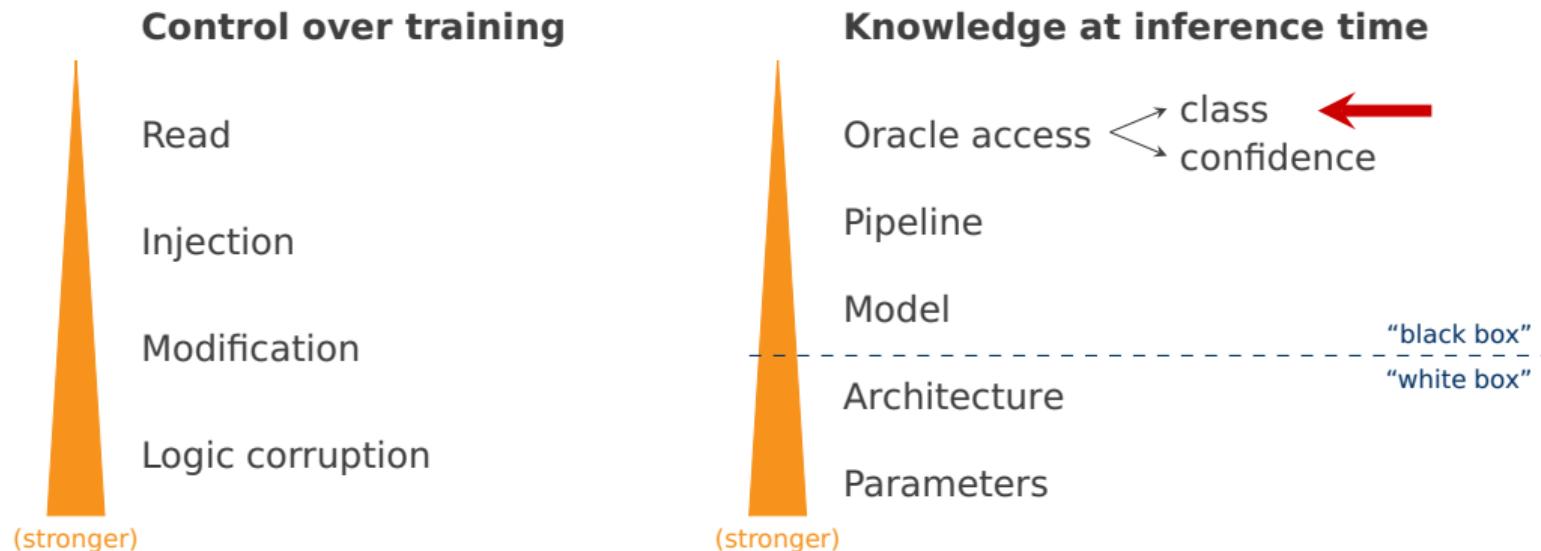
$|\mathbf{w}|$ size of the classifier (e.g., # of perceptrons in an ANN)

d distance to decision boundary

k dimension of the data

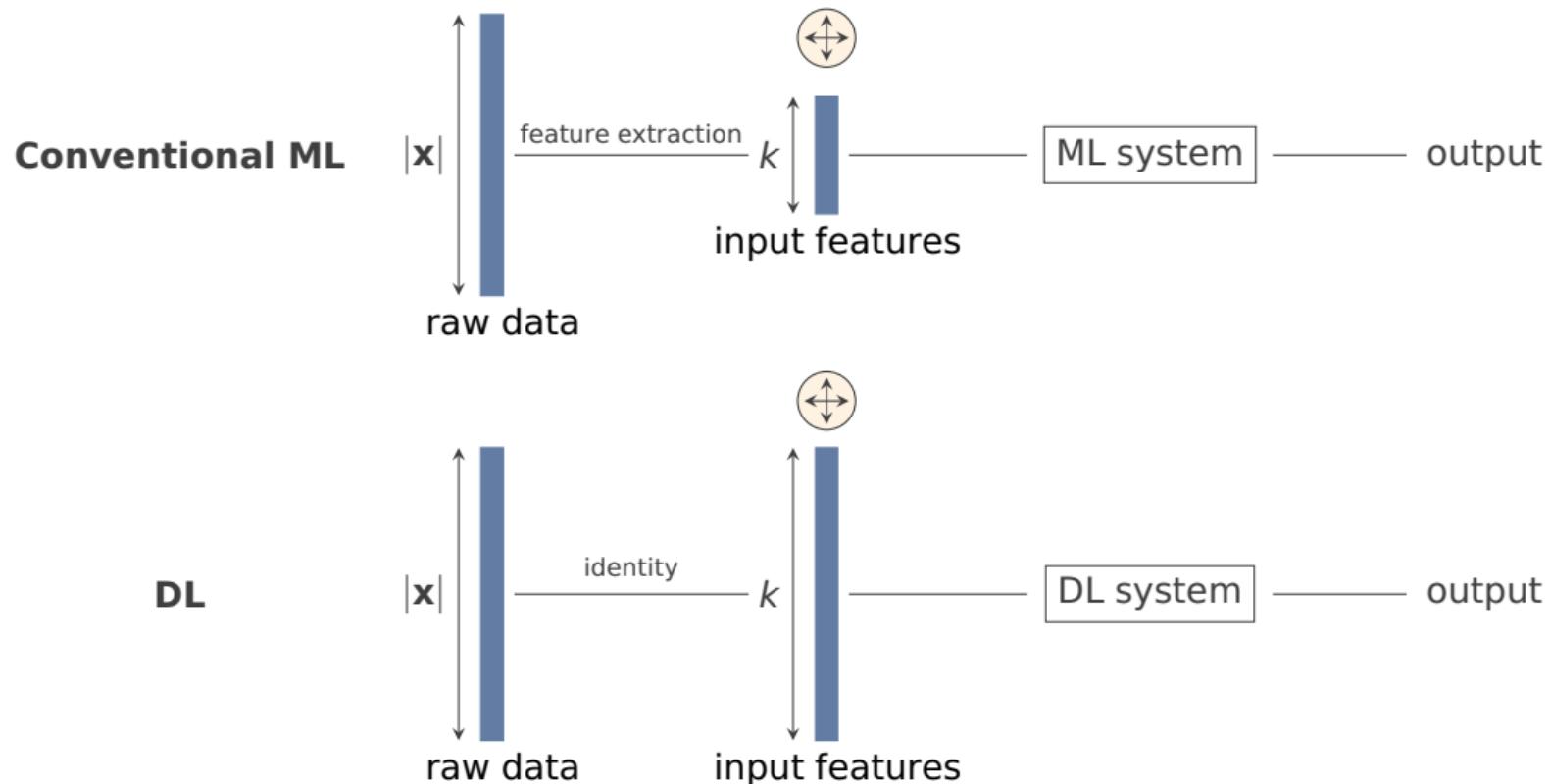
Memo item: all attacks on this slide are fast and approximate. The optimal attack is almost always NP-hard in k .

Adversarial Capabilities



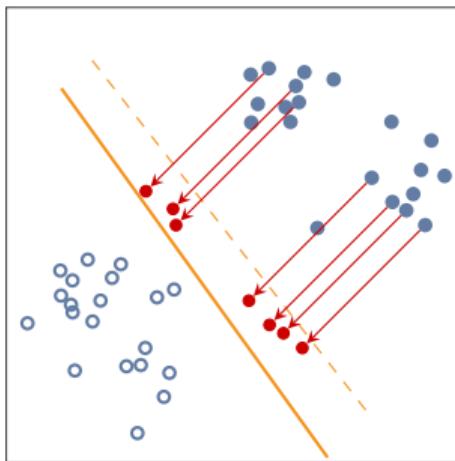
Adapted from Papernot et al. 2018, Fig. 3

Specific Weakness of Deep Learning



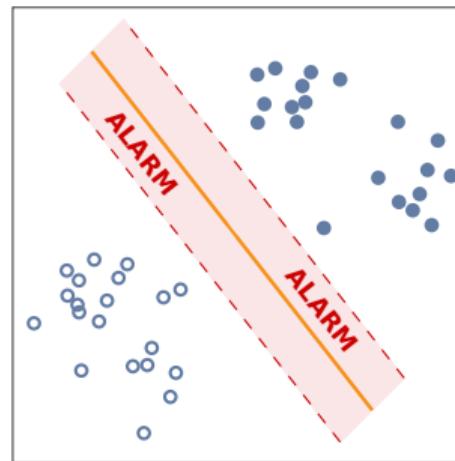
Defenses Against Evasion Attacks

Adversarial training



Attacker's response: adjust scaling

Detection at inference time

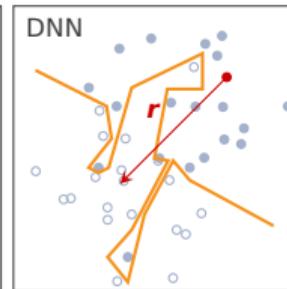
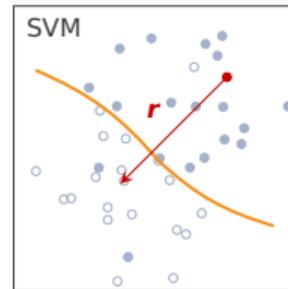
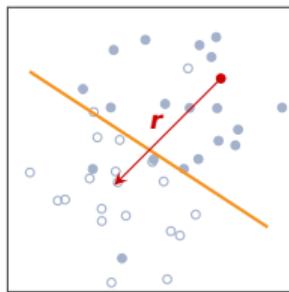


Attacker's response: add alarm as constraint

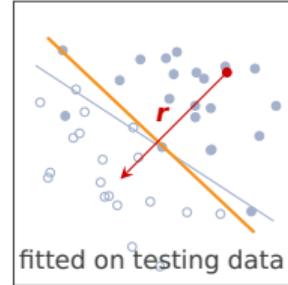
Curated list of white-box defenses against evasion attacks: <https://www.robust-ml.org/defenses/>

Transferability

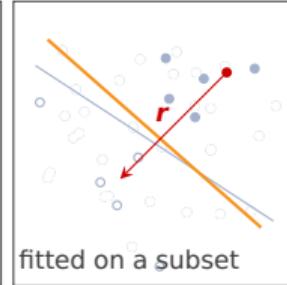
to other classifiers:



to substitute models:

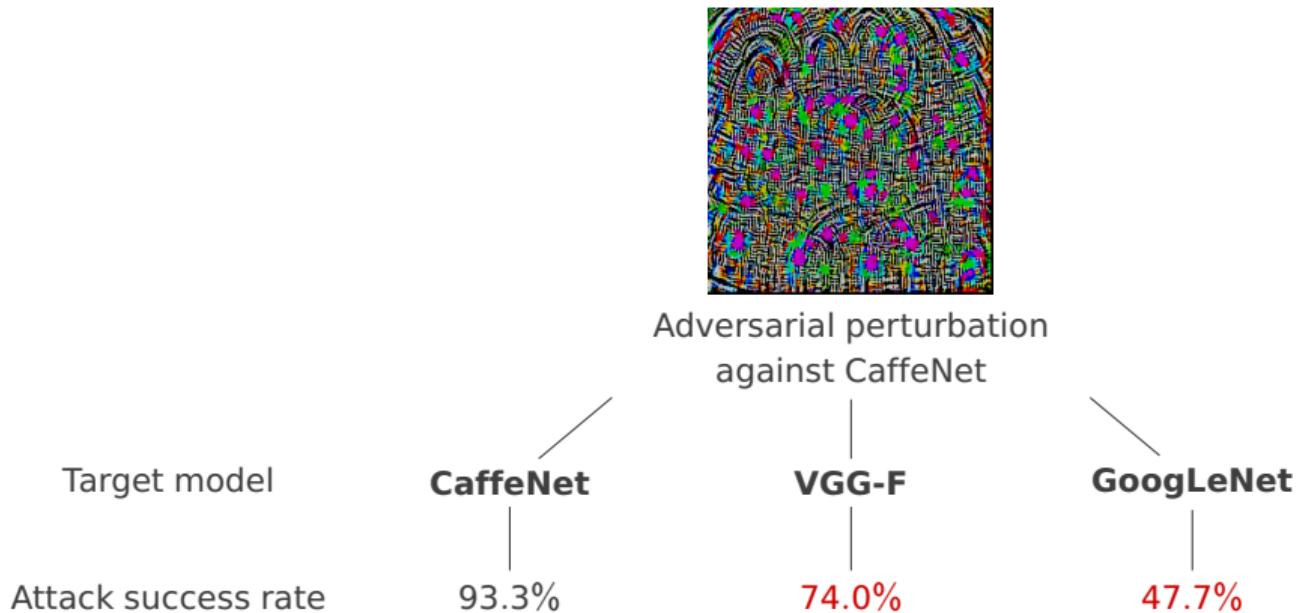


fitted on testing data



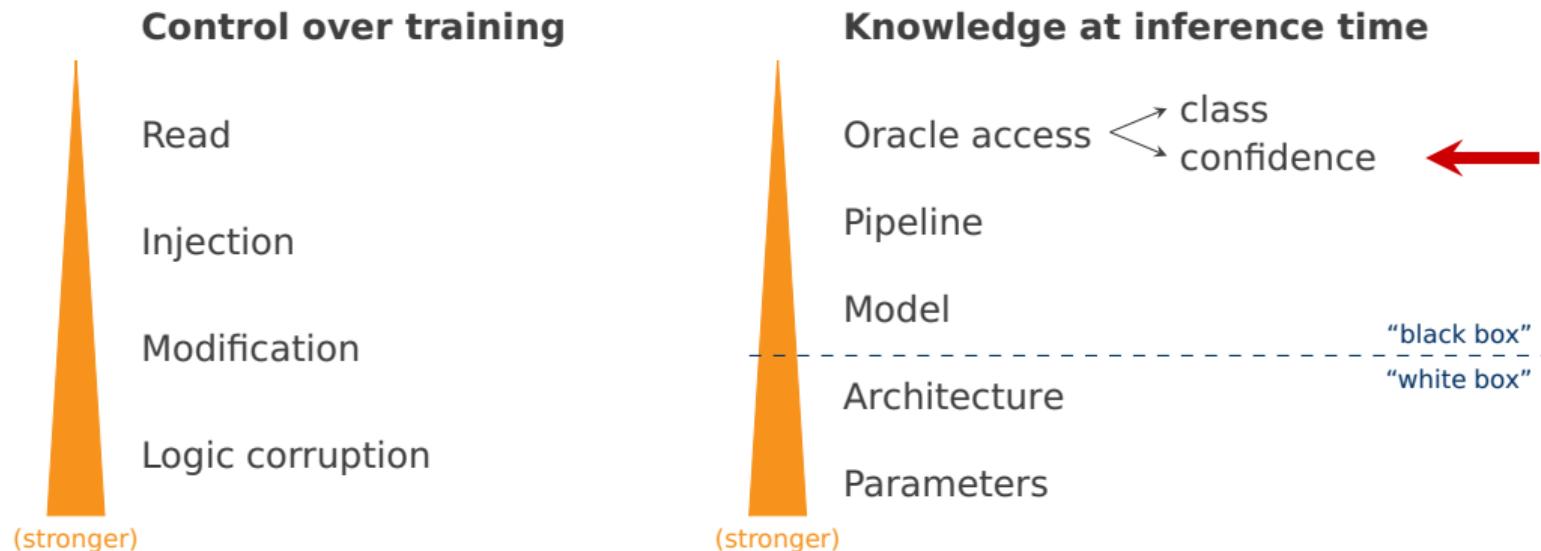
fitted on a subset

Example for Transferability



Moosavi-Dezfooli et al. (2017): Universal adversarial perturbations, *IEEE CVPR*.

Adversarial Capabilities



Adapted from Papernot et al. 2018, Fig. 3

Stages in Supervised Learning

Training

Find model parameters θ that minimize the dissimilarity of model outputs $h_\theta(\mathbf{x})$ and the corresponding known y .

$$\arg \min_{\theta} \left| h_{\theta}(\mathbf{x}_i) - y_i \right|_i$$

Inference

Fix θ . Predict property of interest for new input \mathbf{x} by computing $h_\theta(\mathbf{x})$.

$$\hat{y} = h_{\theta}(\mathbf{x})$$

read access: model stealing

h_{θ}
 θ

machine learning model

parameters

\mathcal{H} candidate hypotheses $\{\mathbf{x} \mapsto h_{\theta}(\mathbf{x}) \mid \theta \in \Theta\}$

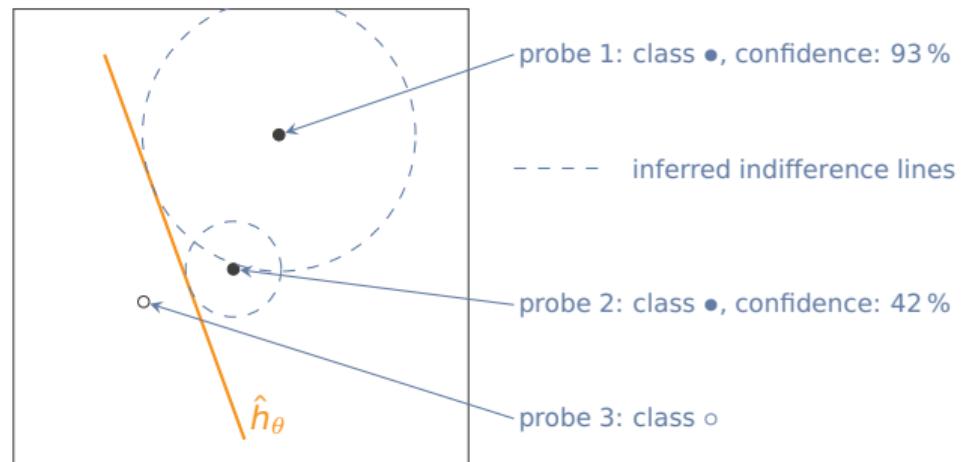
$|\cdot|_i$ loss function on all training points $i = 1, \dots$

\hat{y} prediction

(\mathbf{x}, y) training point: (input, output)

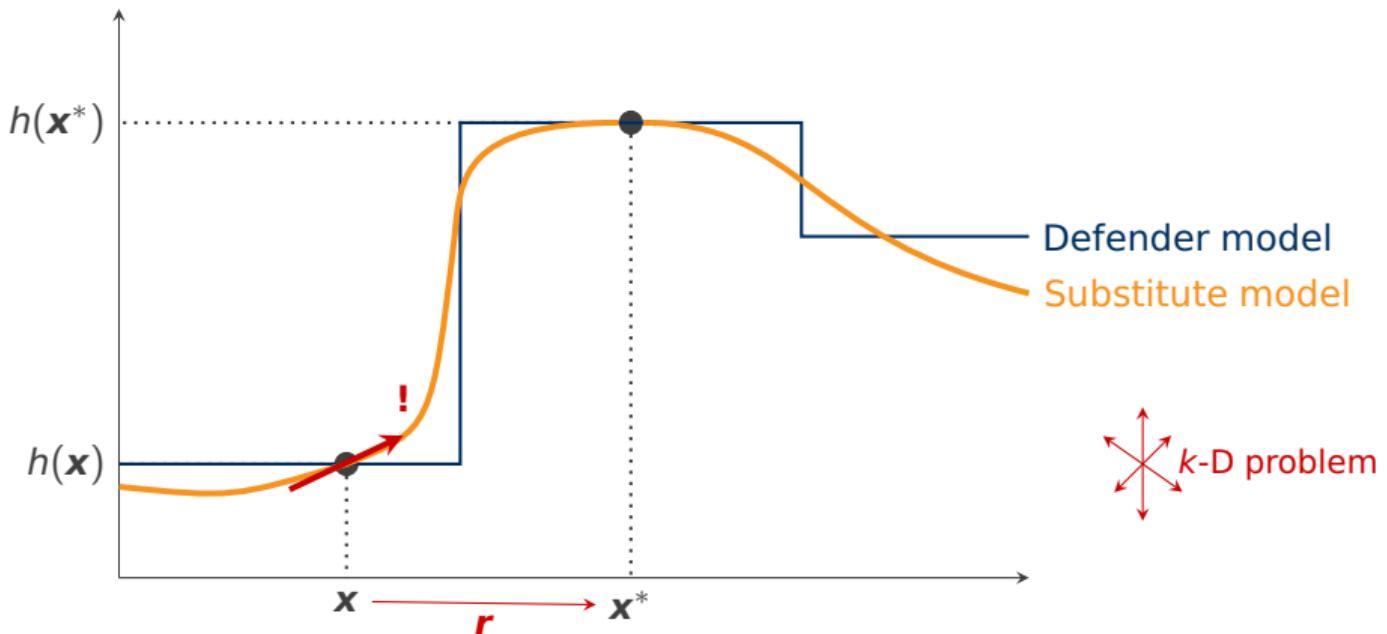
Model Stealing

Example for $k = 2$ dimensions, inference time



General approach: query $O(k)$ linear independent probes and solve equation system

Gradient Masking



Cf. Papernot et al. 2018, Fig. 5

Stages in Supervised Learning

Training

Find model parameters θ that minimize the dissimilarity of model outputs $h_\theta(\mathbf{x})$ and the corresponding known y .

$$\arg \min_{\theta} \left| h_{\theta}(\mathbf{x}_i) - y_i \right|_i$$

Inference

Fix θ . Predict property of interest for new input \mathbf{x} by computing $h_\theta(\mathbf{x})$.

$$\hat{y} = h_{\theta}(\mathbf{x})$$

h_θ machine learning model

θ parameters

\mathcal{H} candidate hypotheses $\{\mathbf{x} \mapsto h_{\theta}(\mathbf{x}) \mid \theta \in \Theta\}$

$\|\cdot\|_i$ loss function on all training points $i = 1, \dots$

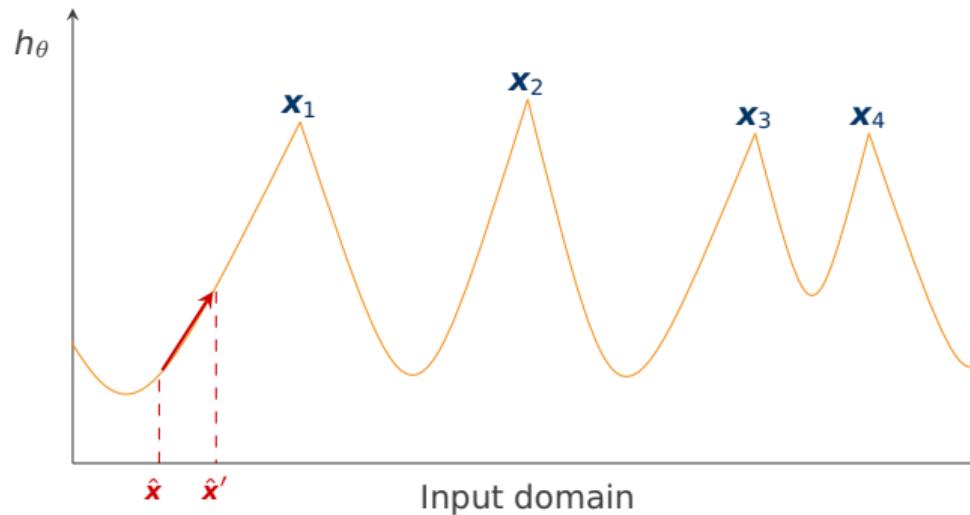
(\mathbf{x}, y) training point: (input, output)

\hat{y} prediction

read access: model inversion

Model Inversion

Training data points are local maxima of h_θ :



Gradient-based methods known from evasion attacks work for model inversion, too.

Example for Model Inversion Attacks

against a face recognition system



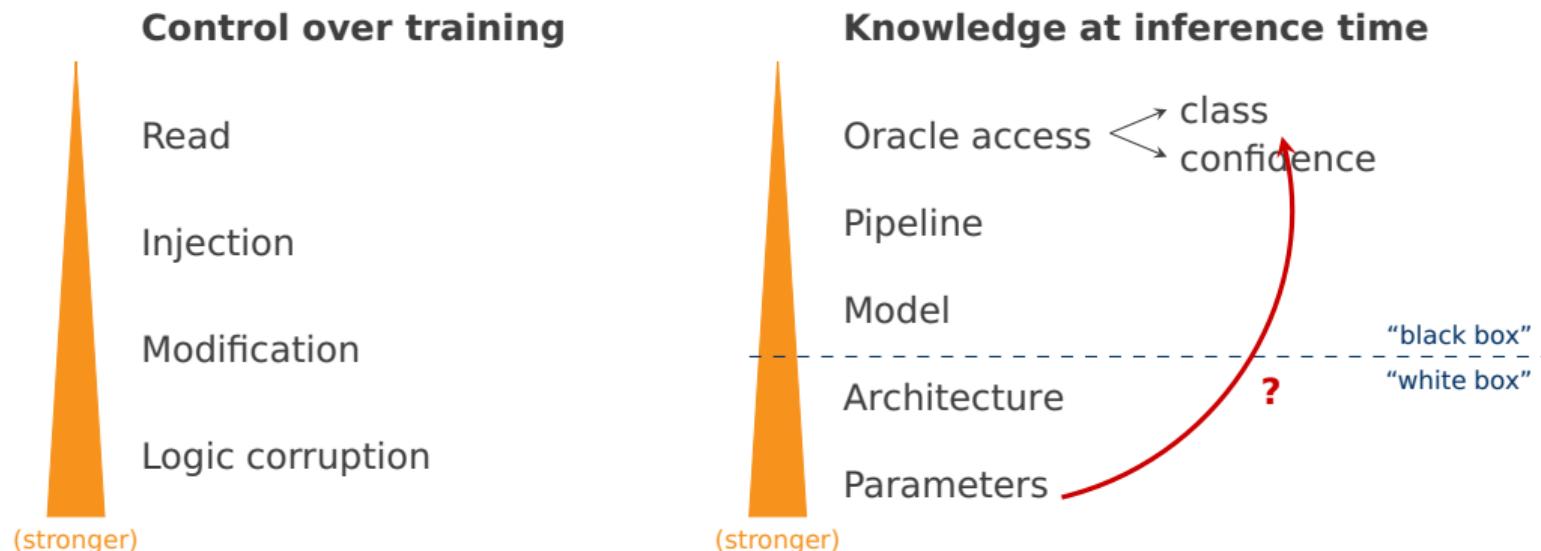
Training point



Reconstructed face

Fredrickson et al. (2015): Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures. ACM CCS.

Adversarial Capabilities



Adapted from Papernot et al. 2018, Fig. 3

Turning White Box Into Black Box

Homomorphic encryption exists:

$$\text{enc}(a) + \text{enc}(b) = \text{enc}(a + b) \quad \text{xor} \quad \text{enc}(a) \cdot \text{enc}(b) = \text{enc}(a \cdot b)$$

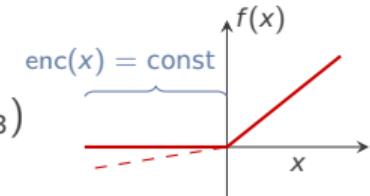
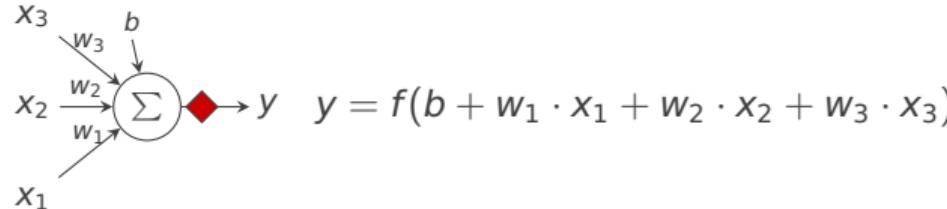
sometimes with efficient ways to calculate:

$$c \cdot \text{enc}(a) = \text{enc}(c \cdot a) \quad \text{xor} \quad \text{enc}(a)^c = \text{enc}(a^c)$$

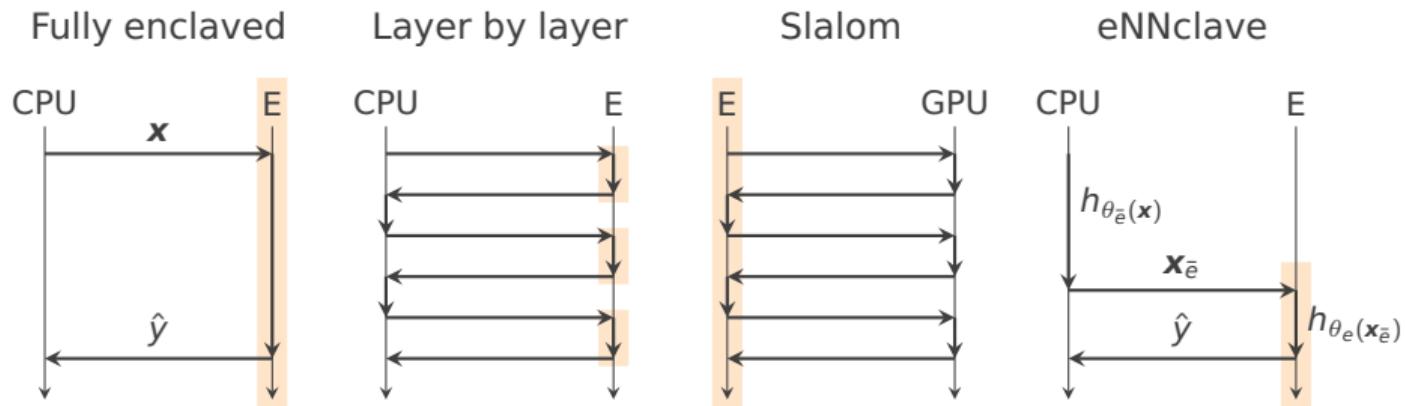
Fully homomorphic encryption is not practical: (1 bit → megabytes of ciphertext)

$$\text{enc}(a) + \text{enc}(b) = \text{enc}(a + b) \quad \text{and} \quad \text{enc}(a) \cdot \text{enc}(b) = \text{enc}(a \cdot b)$$

Now consider:



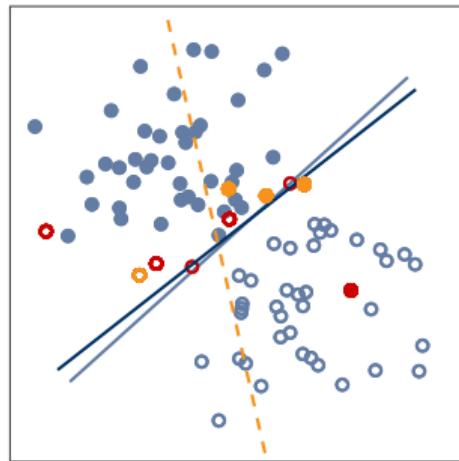
Using Trusted Processors



Ohrimenko et al. 2016, Hanzlik et al. 2019, Tramèr & Boneh 2019, Schlögl & Böhme 2020 (Fig. 6)

Fairness

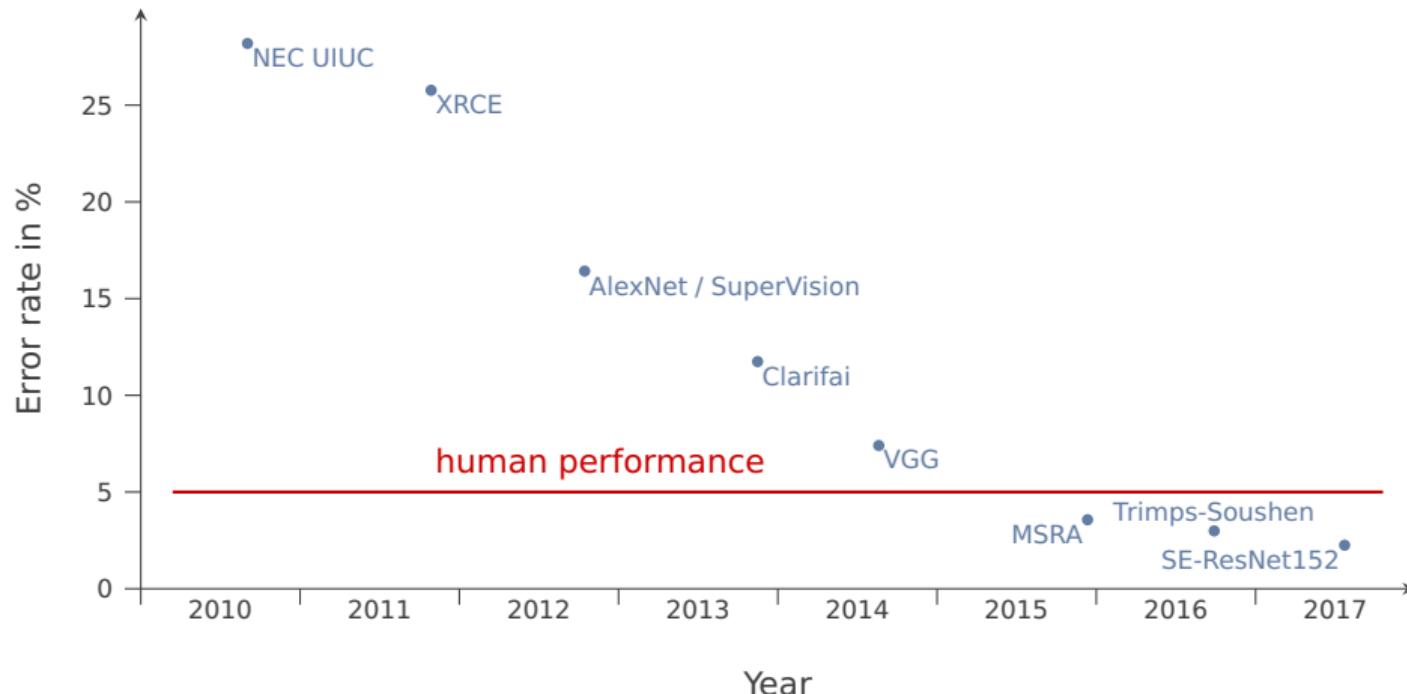
Population



Minority

Trained on	Error rates (in %)	
	Majority	Minority
Majority	2.5	75.0
All	2.5	62.5

Progress in Image Recognition

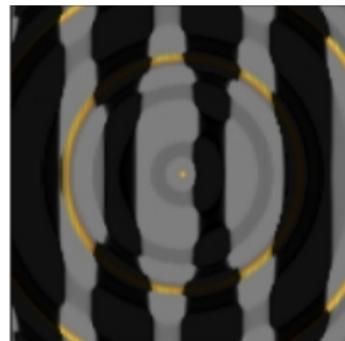


Benchmark: Imagenet

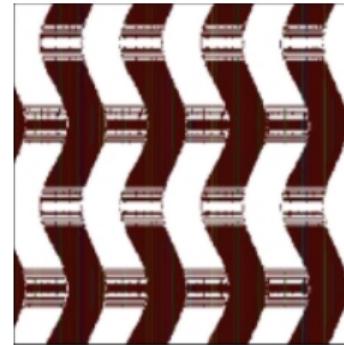
Data source: Electronic Frontier Foundation 2018

Superior to Human Performance ?

Maybe we should rethink our measurement criteria.



“penguin”
confidence: 99.9 %



“electric guitar”
confidence: 99.9 %

A. Nguyen et al. Deep Neural Networks are Easily Fooled. 2015 *IEEE CVPR*.

Lessons Learned

Machine learning can generalize from training samples to inference samples if they are drawn from the same distribution.

An adversary with (partial) control over any one of these distributions may not only compromise performance, but can control outputs (predictions) to a certain extent.

It is hard to confine inputs to a machine learning function with regard to any other use of the same data or model.

The adversary's task seems to become easier the more data-driven the model is, i.e., the fewer assumptions are given from the outside.

Obtaining some (defensible) guarantees for inference-based systems is currently an area of very active research.

Syllabus – Summer Term 2021

- | | |
|-----------------|--|
| 02.03.21 | 1. Introduction, prerequisites, a secure single-purpose device |
| 09.03.21 | 2. Machine learning in adversarial environments |
| 16.03.21 | 3. Multi-purpose systems: confinement & side channels |
| 23.03.21 | 4. Multi-purpose systems: access control & vulnerabilities |
| 13.04.21 | 5. Hardware-supported security systems |
| 20.04.21 | 6. Securing end-to-end network connections |
| 27.04.21 | 7. Securing network infrastructures |
| 04.05.21 | 8. Availability |
| 11.05.21 | 9. Security economics |
| 18.05.21 | 10. Privacy policy & theory |
| 01.06.21 | 11. Privacy-enhancing technology |
| 08.06.21 | 12. Q & A |
| 22.06.21 | Oral exams |