

# Adversarial Machine Learning (AML)

Somesh Jha

University of Wisconsin, Madison

ICISS 2018 (Bangalore)

*Thanks to Nicolas Papernot, Ian Goodfellow, and Jerry Zhu for some slides.*

# Thanks

- Collaborators...
- NSF
  - SaTC Frontiers Grant (Penn State, UCSD, Stanford, Virginia, Berkeley, Wisconsin)
  - <https://ctml.psu.edu/>
  - FMitF

- ARO

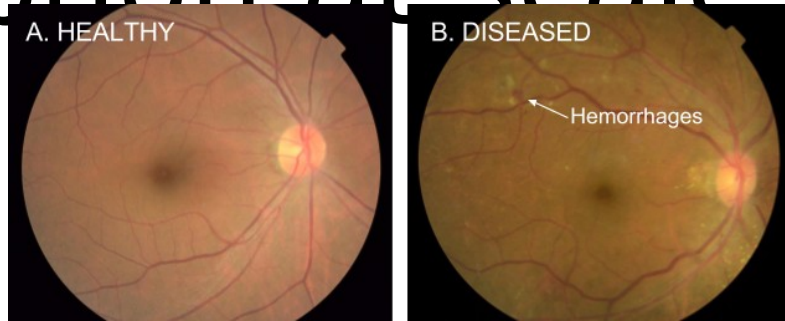


# Announcements/Caveats



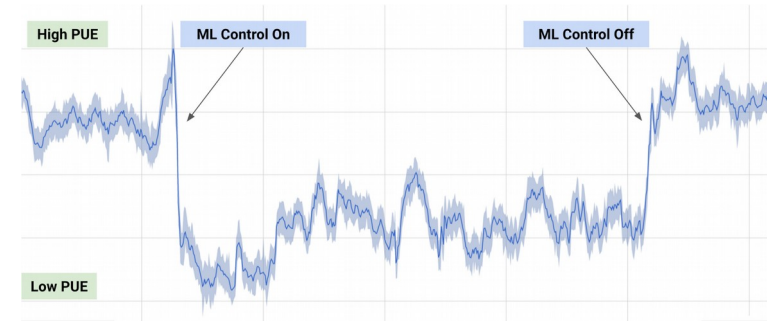
- Please ask questions during the talk
  - If we don't finish, fine☐
- More slides than I can cover
  - Lot of skipping will be going on
- Fast moving area
  - Apologies if I don't mention your paper
- Legend

# Machine learning brings social disruption at scale



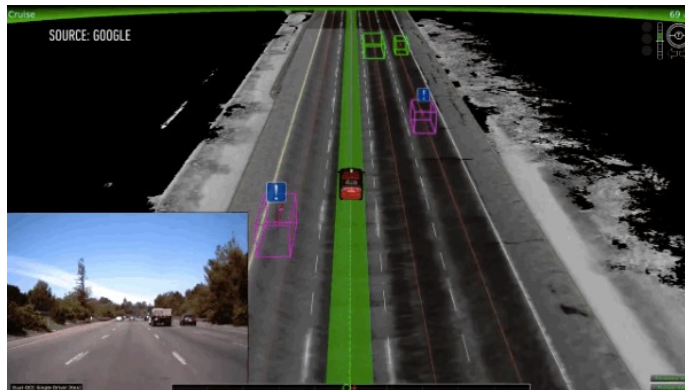
## Healthcare

Source: Peng and Gulshan (2017)



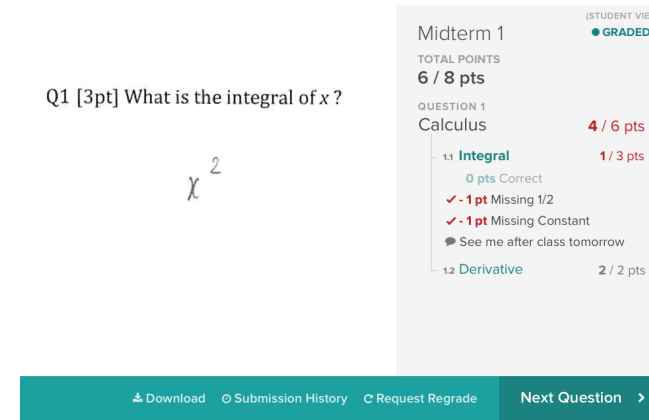
## Energy

Source: Deepmind



## Transportation

Source: Google

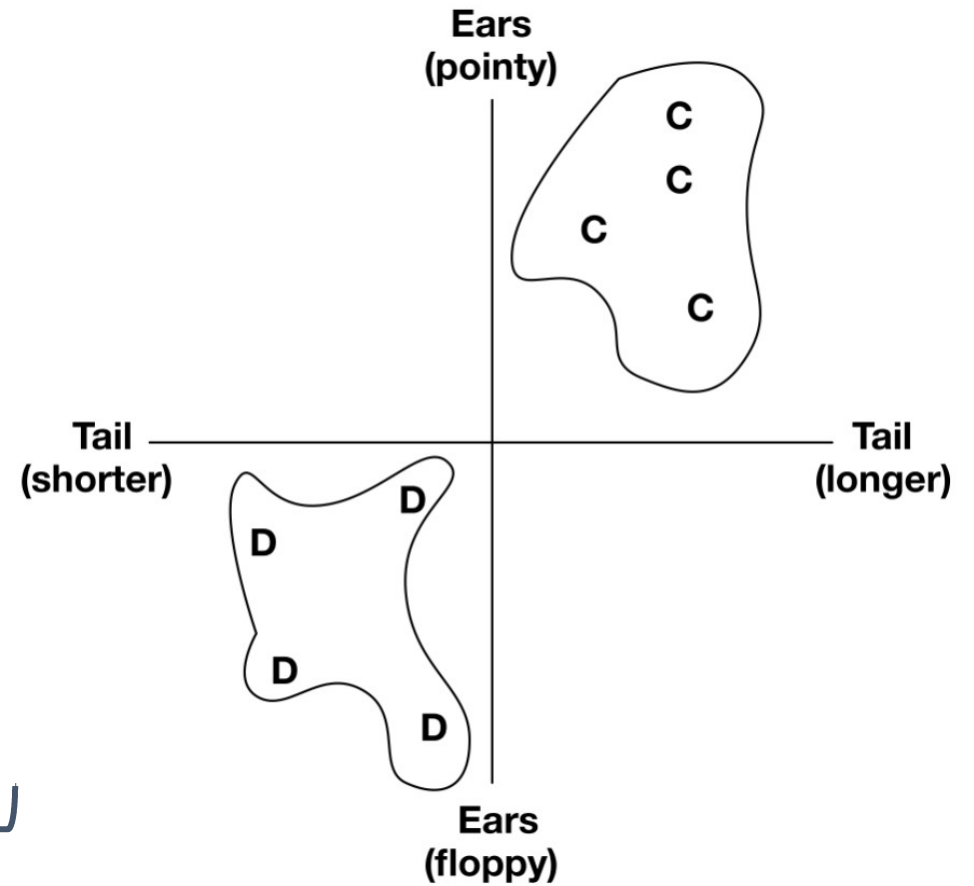


The screenshot shows a Gradescope interface for a 'Midterm 1'. The question is 'Q1 [3pt] What is the integral of  $x^2$ ?'. The student's answer is  $x^2$ . The interface shows the question is graded, with a total score of 6/8 points. The question score is 4/6 points. The student's answer is marked as incorrect (0 pts correct) with a note: '-1 pt Missing 1/2' and '-1 pt Missing Constant'. There is a note 'See me after class tomorrow'. The next question is 'Derivative' worth 2/2 points. At the bottom, there are buttons for 'Download', 'Submission History', 'Request Regrade', and 'Next Question'.

## Education

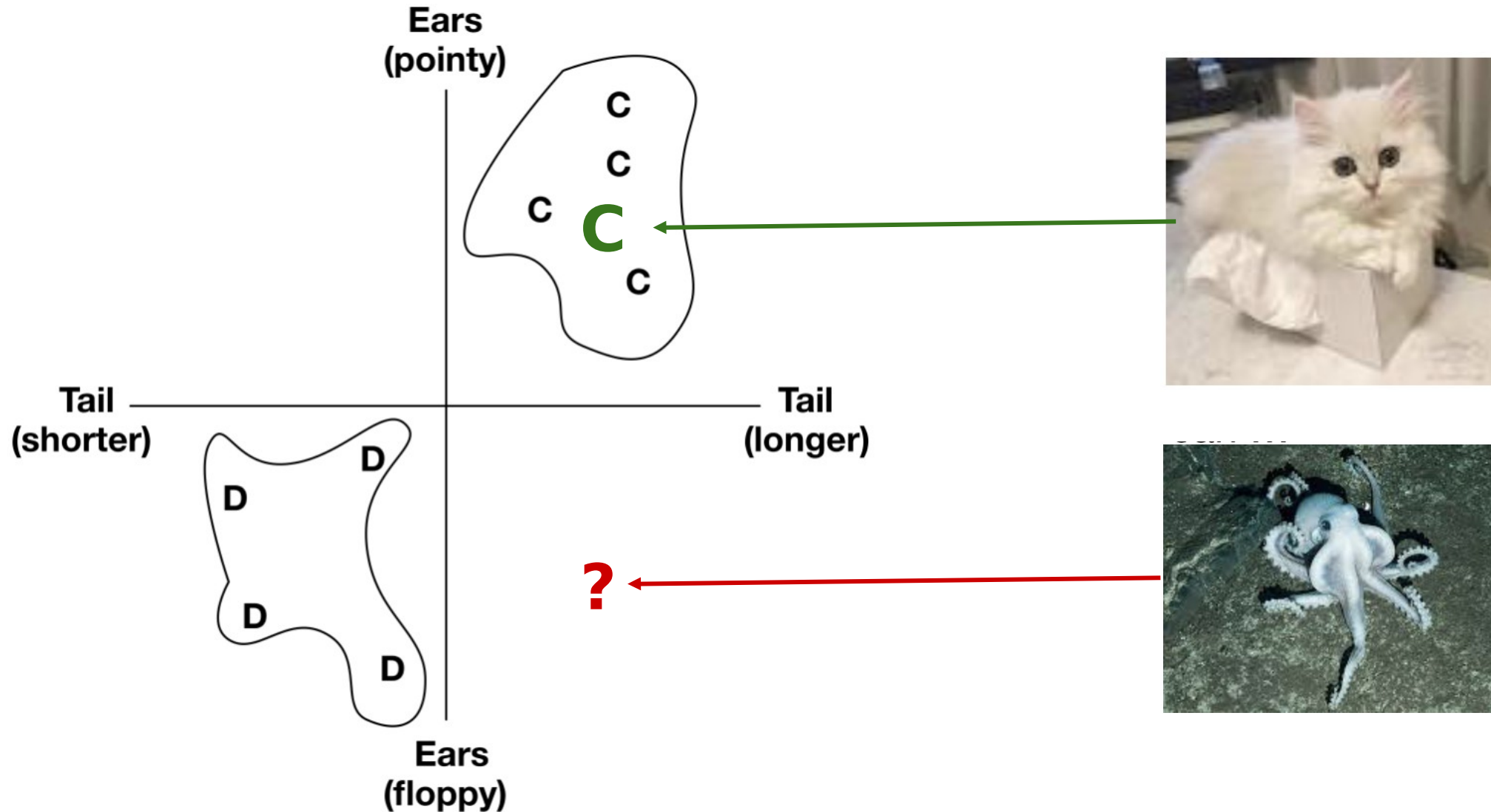
Source: Gradescope

# Machine learning is not magic (training time)

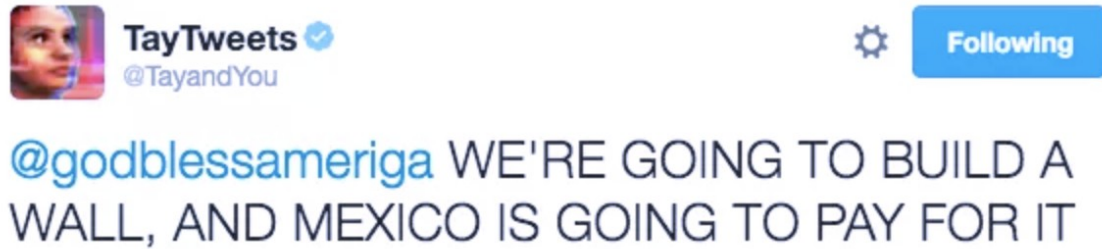


**Training data**

# Machine learning is not magic (inference time)



# Machine learning is deployed in adversarial settings



## Microsoft's Tay chatbot

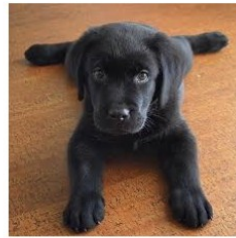
*Training* data poisoning



## YouTube filtering

Content evades detection at *inference*

# Machine learning does not always generalize well

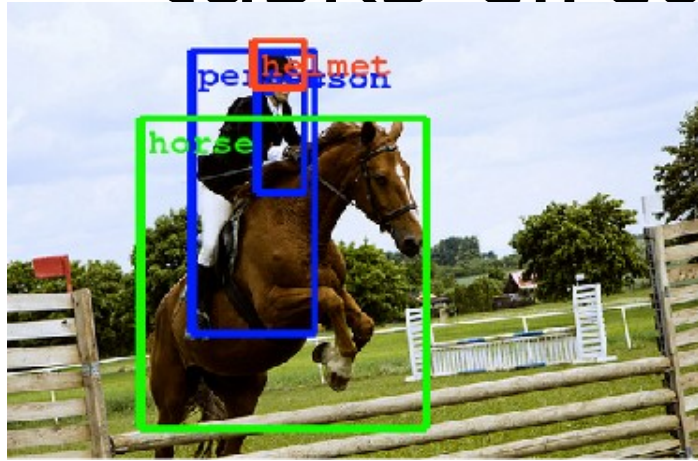


**Training data**

**Test data**



# ML reached “human-level performance” on many IID tasks circa 2013



(Szegedy et al, 2014)

...recognizing objects and faces....

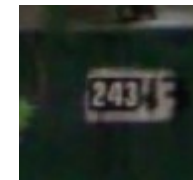


(Taigmen et al, 2013)



(Goodfellow et al, 2013)

...solving CAPTCHAS and reading addresses...

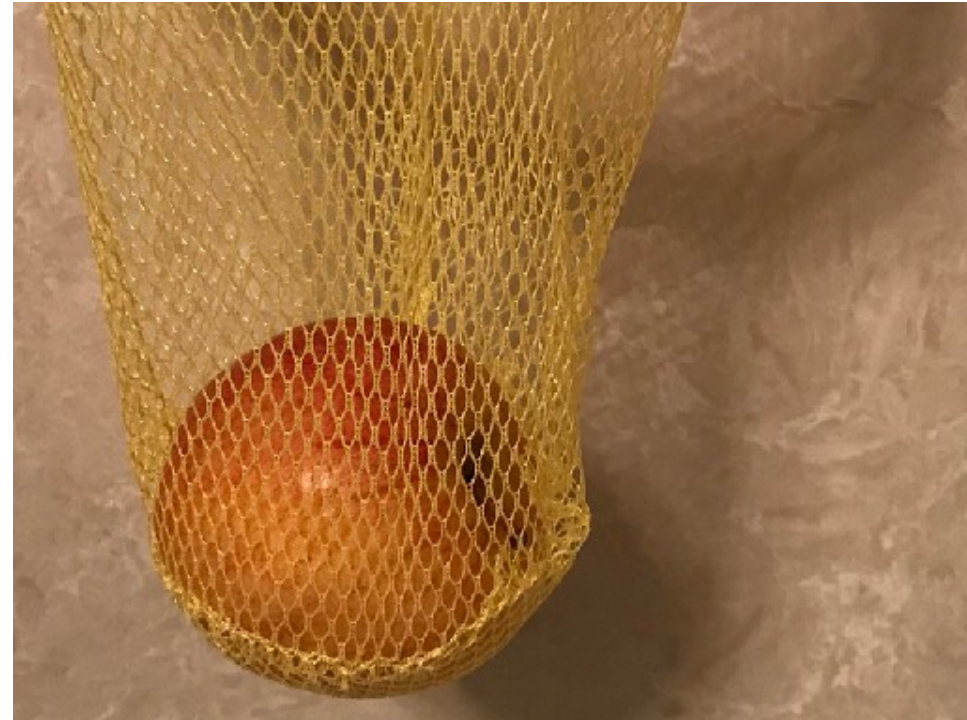


(Goodfellow et al, 2013)

# Caveats to “human-level” benchmarks



Humans are not very good at some parts of the benchmark



The test data is not very diverse. ML models are fooled by natural but unusual data.

# Deluge of Work...

- Help, I can't keep up □
- Attacks
- Defenses
  - Adhoc and certified
- Other domains
  - Text, malware, ....
- Verification algorithm
- ....



# ML (Basics)

- Supervised learning
- Entities
  - (Sample Space)  $Z = X \times Y$ 
    - (data, label)  $(x, y)$
  - (Distribution over  $Z$ )
  - (Hypothesis Space)
  - (Loss function)  $(H \times Z) \rightarrow R$

# ML (Basics)

- *Learner's problem*

- Find  $w$  that minimizes

- (Regularizer)

- $E_{\{z \sim D\}} l(w, z) + \lambda R(w)$

- $\frac{1}{m} \sum_{\{i=1\}}^m l(w, (x_i, y_i)) + \lambda R(w)$

- **Sample set** =  $\{(x_1, y_1), \dots, (x_m, y_m)\}$

- **SGD**

- (iteration)  $w[t + 1] = w[t] - \eta_t l'(w[t], (x_{\{i_t\}}, y_{\{i_t\}}))$

- (learning rate)  $\eta_t$

- ...

# ML (Basics)

- SGD
  - How learning rates change?
  - In what order you process the data?
    - Sample-SGD
    - Random-SGD
  - Do you process in mini batches?
  - When do you stop?

# ML (Basics)

- After Training

- $F_w: X \rightarrow Y$

- $F_w(x) = \operatorname{argmax}_{\{y \in Y\}} s(F_w)(x)$

- (softmax layer)  $(F_w)$

- Sometimes we will write  $F$  simply as  $F$  as

- will be implicit

# ML (Basics)

## • Logistic Regression

- $X = \mathcal{R}^n, Y = \{+1, -1\}$

- $H = \mathcal{R}^n$

- Loss function  $l(w, (x, y))$

- $\log(1 + \exp(-y (w^T x)))$

- $R(w) = \|w\|_2$

- Two probabilities  $s(F) = (p_{\{-1\}}, p_{\{+1\}})$

- $\left( \frac{1}{1 + \exp(w^T x)}, \frac{1}{1 + \exp(-w^T x)} \right)$

## • Classification

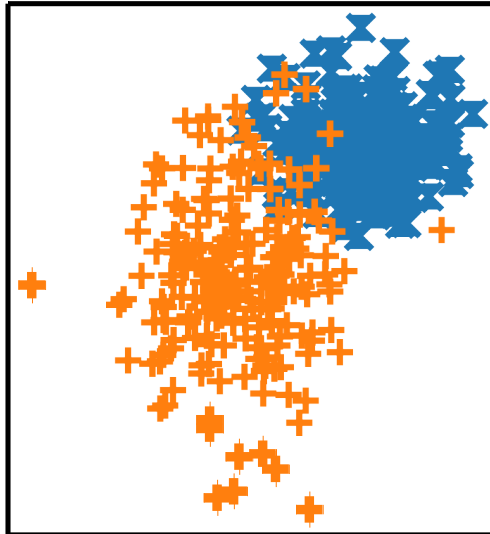
- Predict -1 if  $p_{\{-1\}} > 0.5$

- Otherwise predict +1

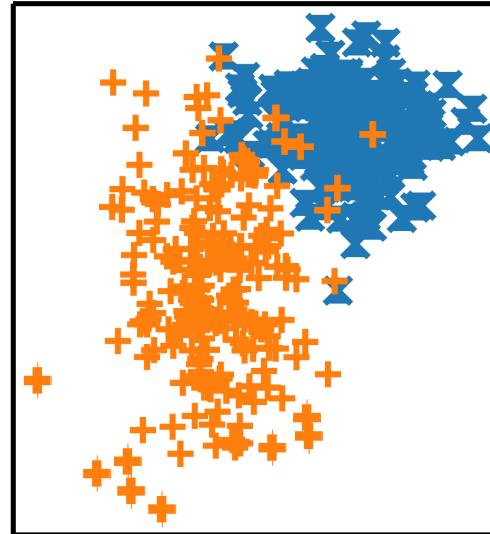


# I.I.D. Machine Learning

Train



Test



I:  
Independent  
I:  
Identically  
D: train and test  
examples drawn  
Distributed  
independently from  
same distribution

# Security Requires Moving Beyond I.I.D.

- Not identical: attackers can use unusual inputs



(Eykholt et al, 2017)

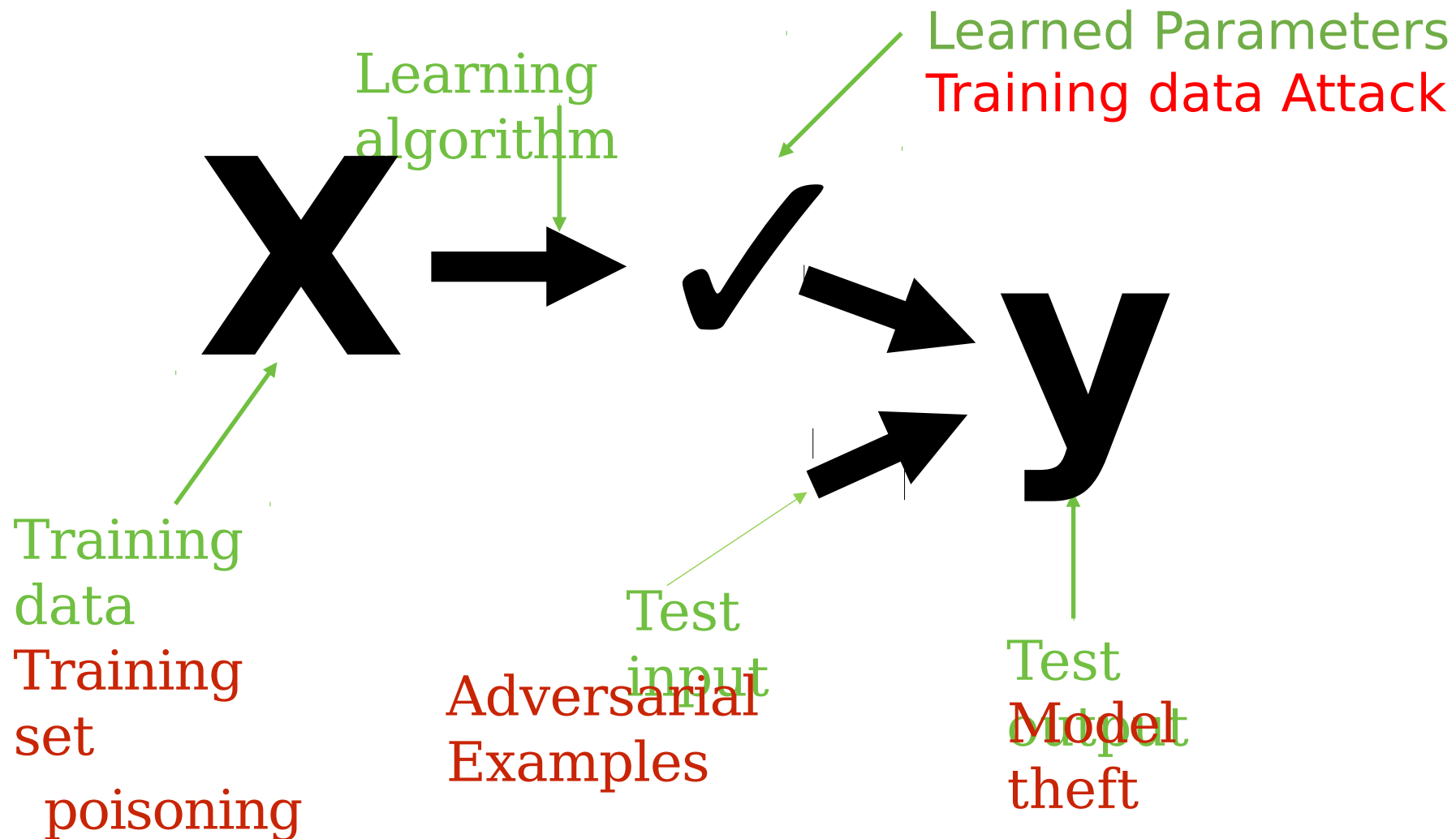
- Not independent: attacker can repeatedly send a single mistake (“test set attack”)

# Adversarial Learning is not new!!

- **Lowd:** I spent the summer of 2004 at Microsoft Research working with Chris Meek on the problem of spam.
  - We looked at a common technique spammers use to defeat filters: adding "good words" to their emails.
  - We developed techniques for evaluating the robustness of spam filters, as well as a theoretical framework for the general problem of learning to defeat a classifier (*Lowd and Meek 2005*)
- But...
  - New resurgence in ML and hence new problems
  - Lot of new theoretical techniques being developed
    - High dimensional robust statistics, robust optimization



# Attacks on the machine learning pipeline



# Fake-News Attacks



# Fake News Attacks

Abusive use of machine learning:

Using GANs to generate **fake content**  
(a.k.a deep fakes)

*Strong societal implications:*

elections, automated trolling, court  
evidence ...

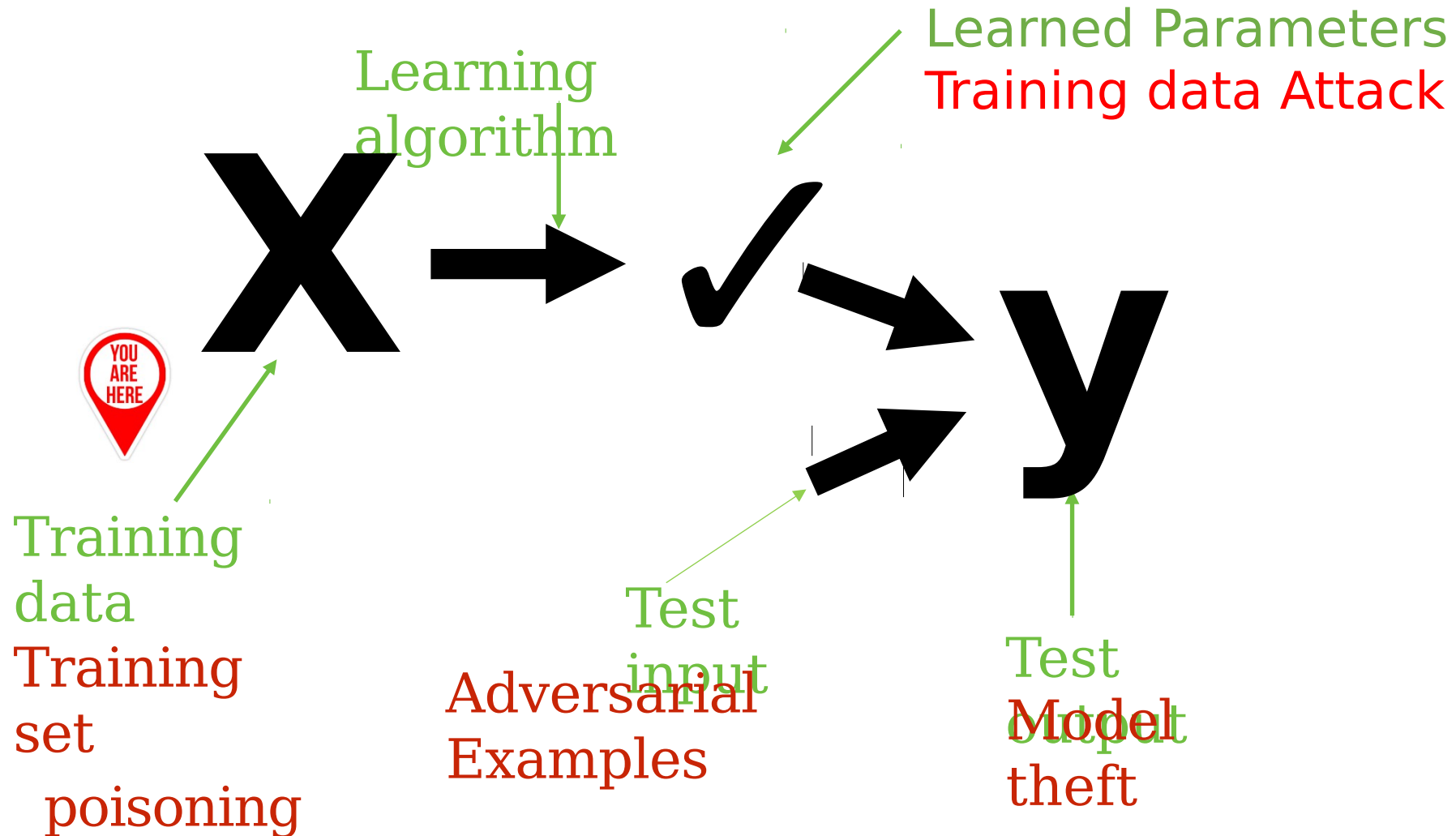


## **Generative media:**

- Video of Obama saying things he never said, ...
- Automated reviews, tweets, comments, indistinguishable from human-generated content

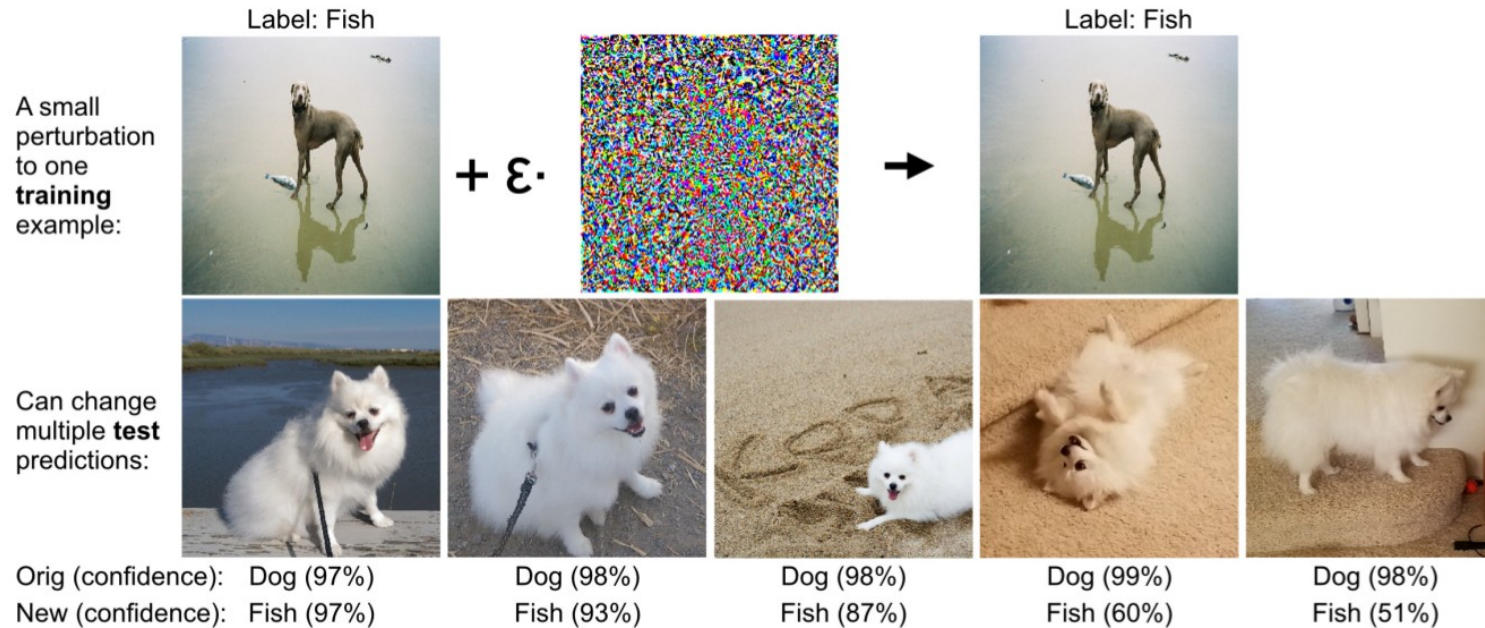
# Training Time Attack

# Attacks on the machine learning pipeline



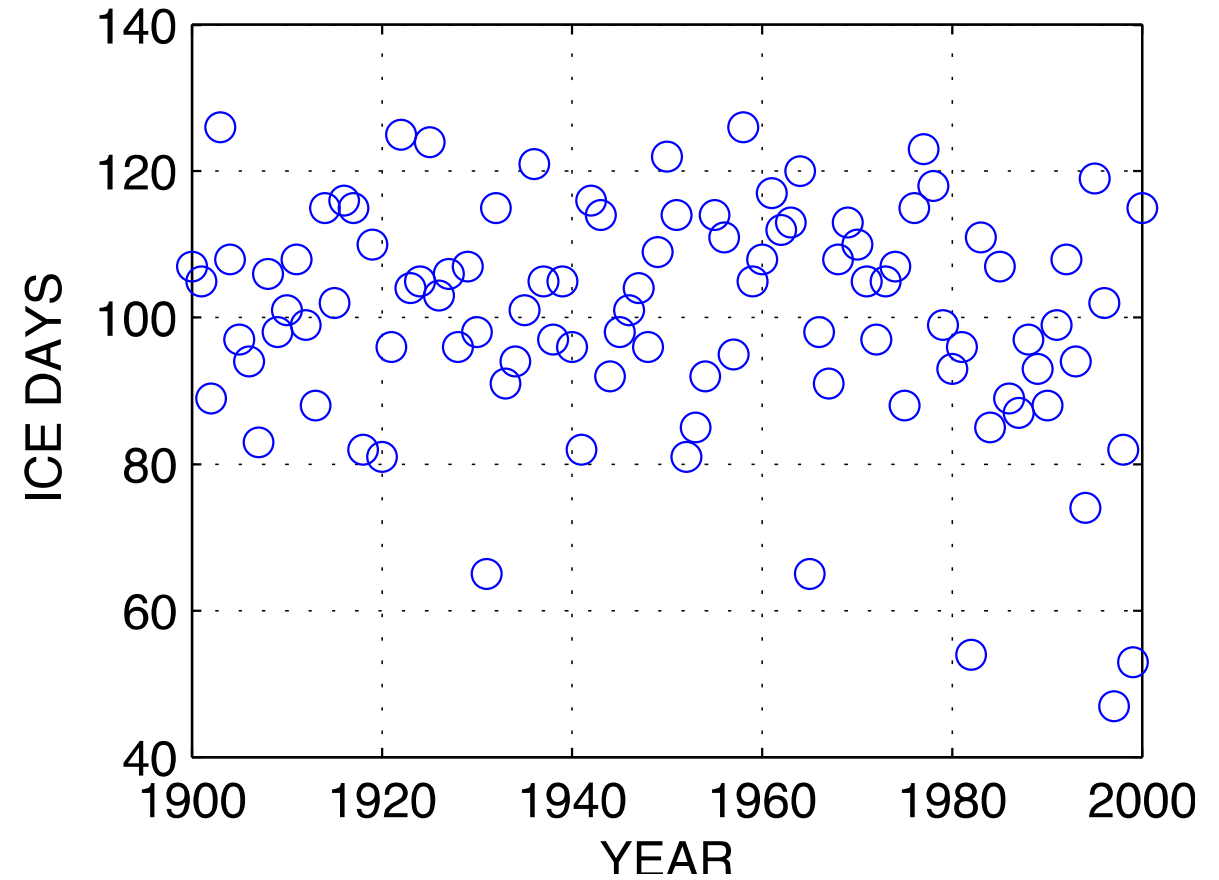


# Training time

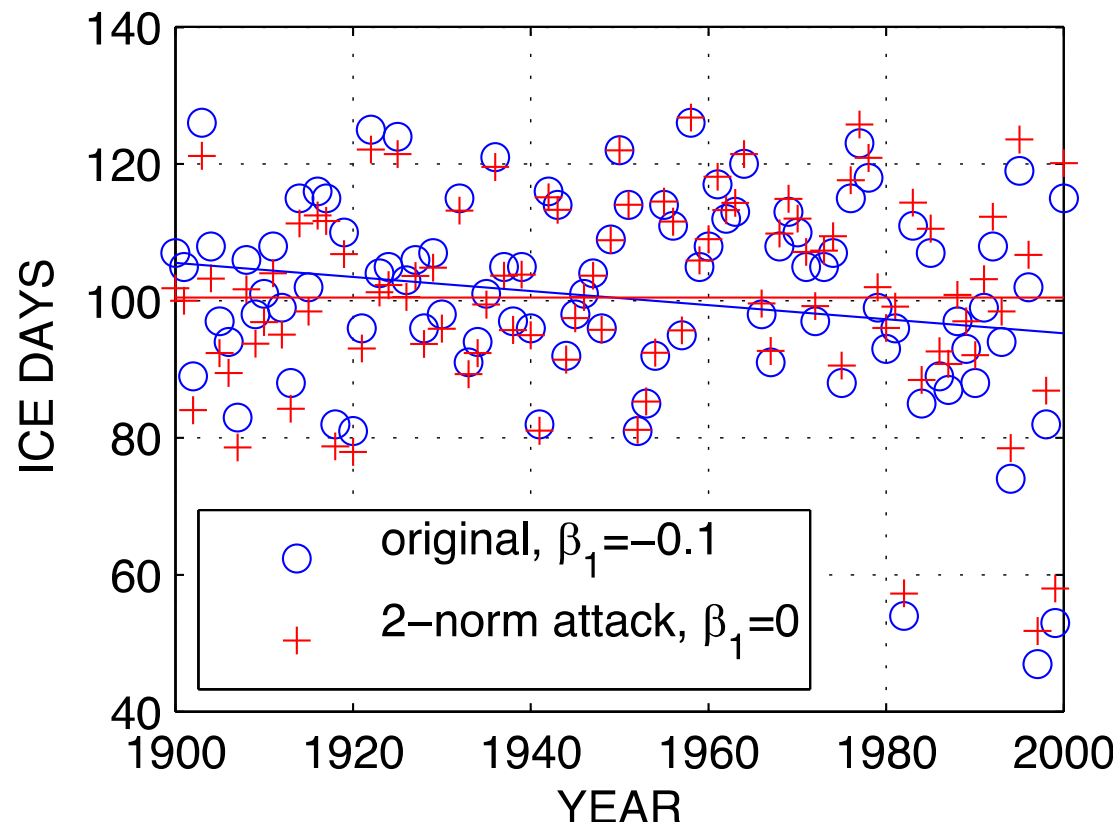
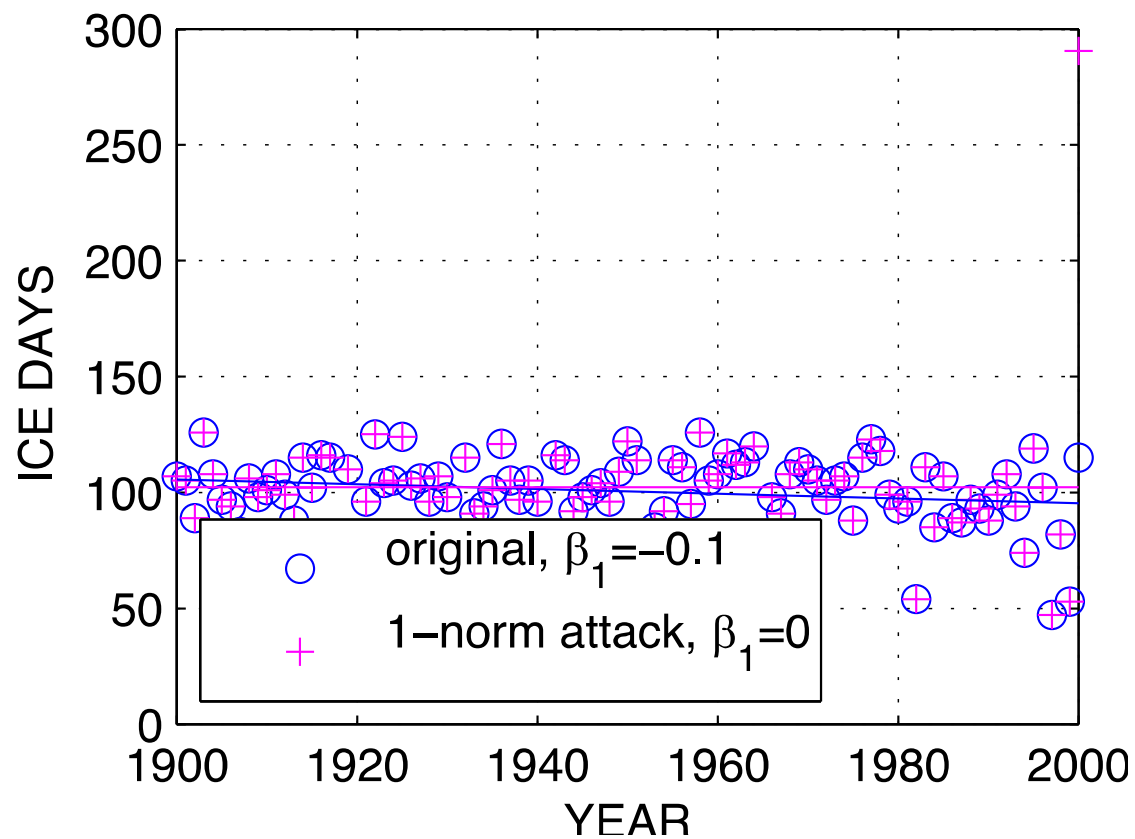


- Setting: attacker perturbs **training set** to fool a model on a test set
- Training data from users is fundamentally a huge security hole
- More subtle and potentially more **pernicious** than test time attacks, due to coordination of multiple points

# Lake Mendota Ice Days



# Poisoning Attacks

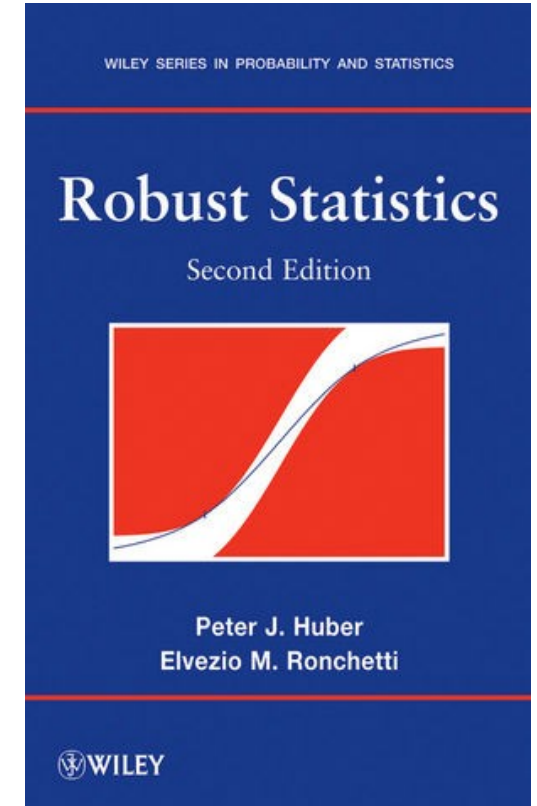


# Formalization

- Alice picks a data set of size  $m$
- Alice gives the data set to Bob
- Bob picks
  - Pick points  $S^B$
  - Gives the data set back to Alice
  - Or could replace some points in  $S$  (less realistic)
- Goal of Bob
  - $S$  is better...
- Goal of Alice
  - Get close to learning from clean data
  - Get close to learning from clean data

# Goal of Bob (bad guy!)

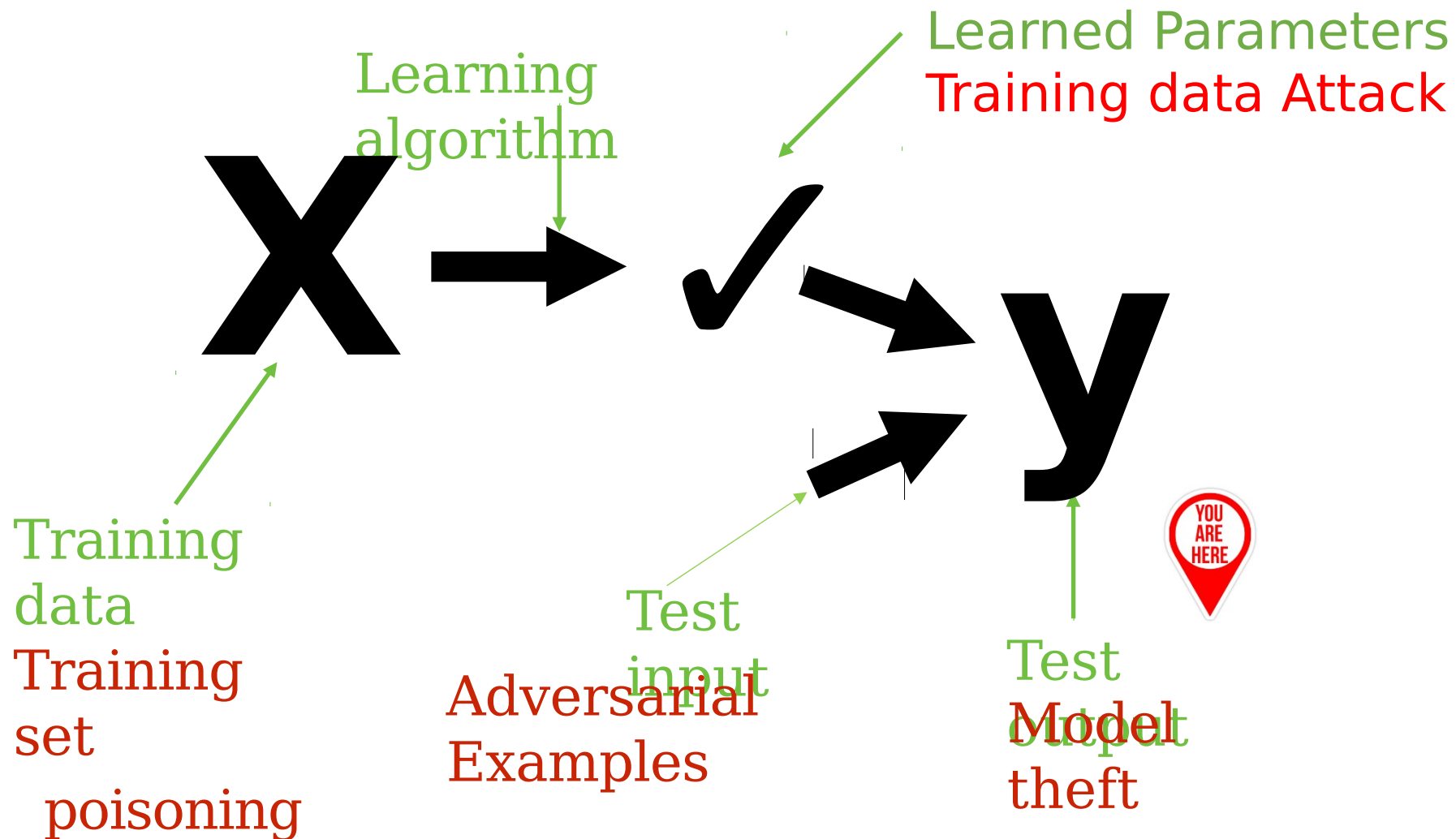
- Maximize the expected value of the loss function
  - Maximize the expected value of the loss function
    - Recall that Alice wants minimize the expected value of the loss function
    - Recall that Alice wants minimize the expected value of the loss function
- Targeted attacks
  - Targeted attacks
    - Picture of *Trent* gets classified as *Vinod*
    - Picture of *Trent* gets classified as *Vinod*
- High-dimensional robust statistics
  - High-dimensional robust statistics
    - $|S^B| = \epsilon m$
    - *Guarantee:* Learn hypothesis that is not “too far” from what you would learn from clean data
    - *Guarantee:* Learn hypothesis that is not “too far” from what you would learn from clean data  $S$



# Representative Papers

- Robust statistics
  - Being Robust (in High Dimensions) Can be Practical  
I. Diakonikolas, G. Kamath, D. Kane, J. Li, A. Moitra, A. Stewart  
ICML 2017
- Certified defenses
  - Certified Defenses for Data Poisoning Attacks. Jacob Steinhardt, Pang Wei Koh, Percy Liang. NIPS 2017
- Targeted attacks
  - Poison Frogs! Targeted Clean-Label Poisoning Attacks on Neural Networks, Shahfi et al., NIPS 2018

# Attacks on the machine learning pipeline



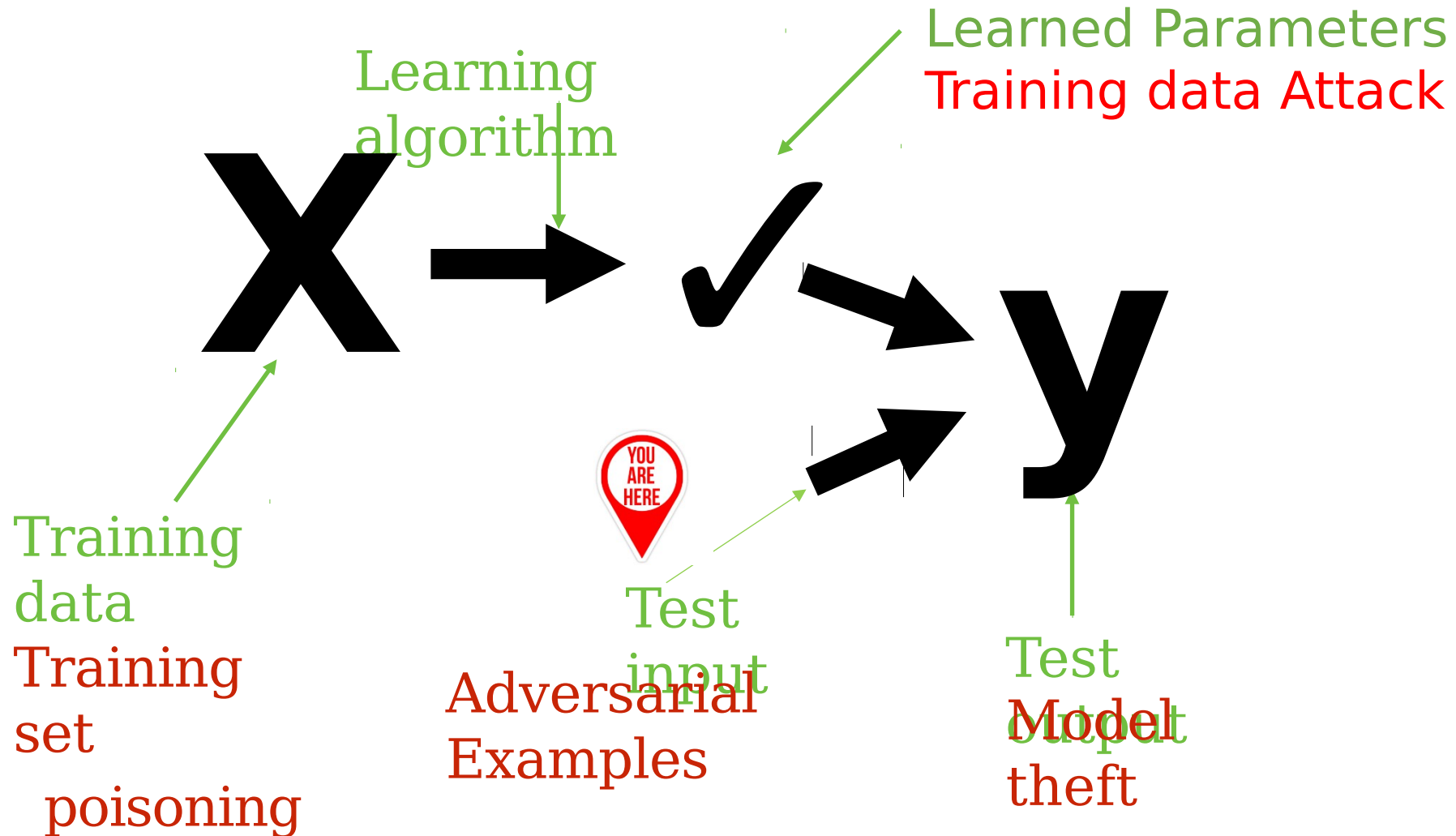
# Model Extraction/Theft Attack



# Model Theft

- **Model theft:** extract model parameters by queries (intellectual property theft)
  - Given a classifier  $F$
  - Query  $F$  on  $q_1, \dots, q_n$  and learn a classifier  $G$
  - $F \approx G$
- **Goals:** leverage active learning literature to develop new attacks and preventive techniques
- **Papers**
  - *Stealing Machine Learning Models using Prediction APIs*, Tramer et al., Usenix Security 2016
  - *Stealing Machine Learning Models using Prediction APIs*, Tramer et al., Usenix Security 2016
  - *Model Extraction and Active Learning*, Chandrasekaran et al.
  - *Model Extraction and Active Learning*, Chandrasekaran et al.

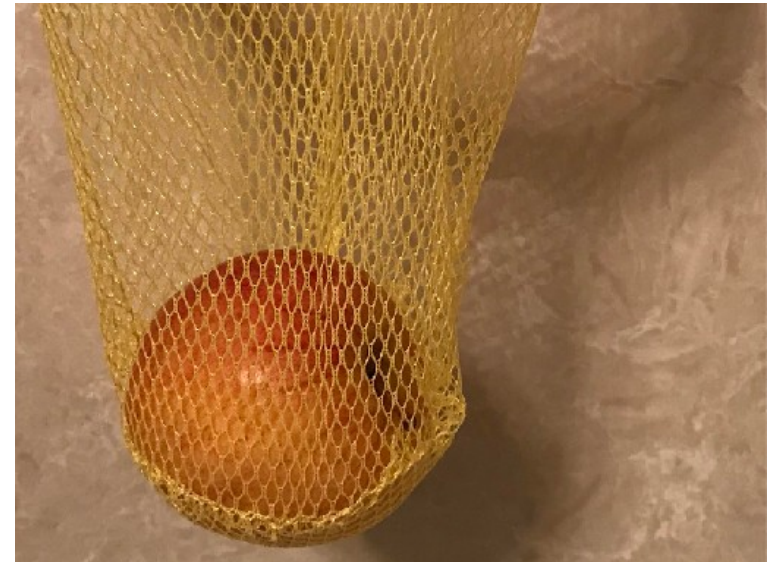
# Attacks on the machine learning pipeline



# Definition

“Adversarial examples are inputs to machine learning models that an attacker has intentionally designed to cause the model to make a mistake”

(Goodfellow et al 2017)



# What if the adversary systematically found these inputs?



$x$   
“panda”  
57.7% confidence

+ .007 ×

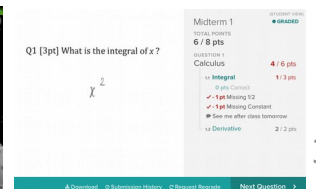
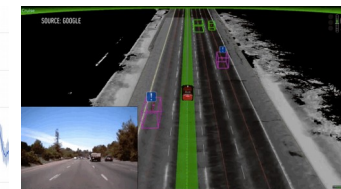
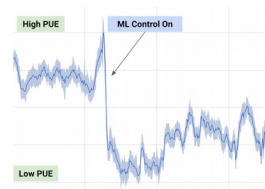
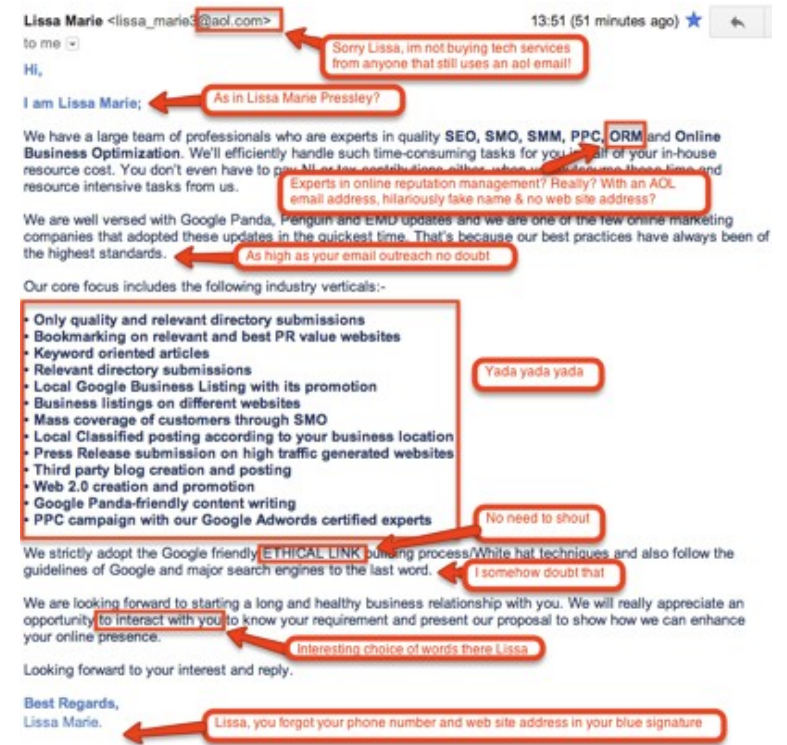


$\text{sign}(\nabla_x J(\theta, x, y))$   
“nematode”  
8.2% confidence

=



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



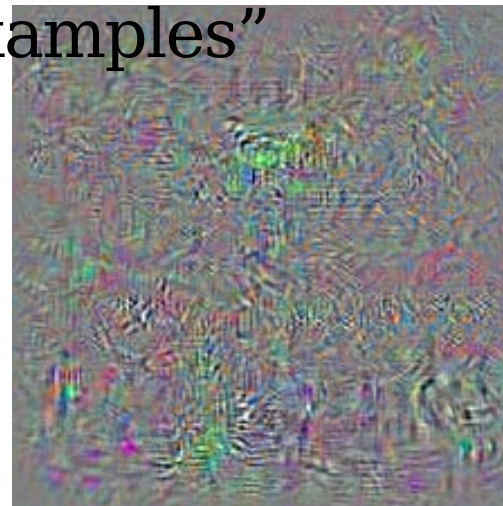
# Good models make surprising mistakes in non-IID setting

“Adversarial  
examples”



Schoolb  
us

+



Perturbation  
(rescaled for visualization)  
(Szegedy et al,  
2013)

=



Ostric  
h

# Adversarial Examples

[Black-box Adversarial Attacks with Limited Queries and Information](#)

Andrew Ilyas, Logan Engstrom, **Anish Athalye**, and Jessy Lin, *ICML 2018*



88% **tabby**  
**cat**

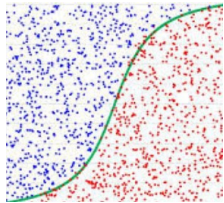


99%  
**guacamole**

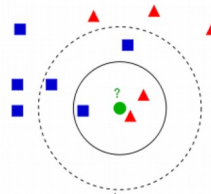
Nice Use of Gradient-Free Optimization

# Adversarial examples...

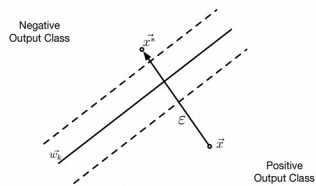
... beyond deep learning



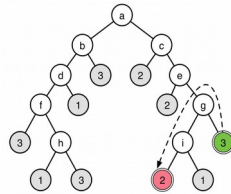
Logistic Regression



Nearest Neighbors

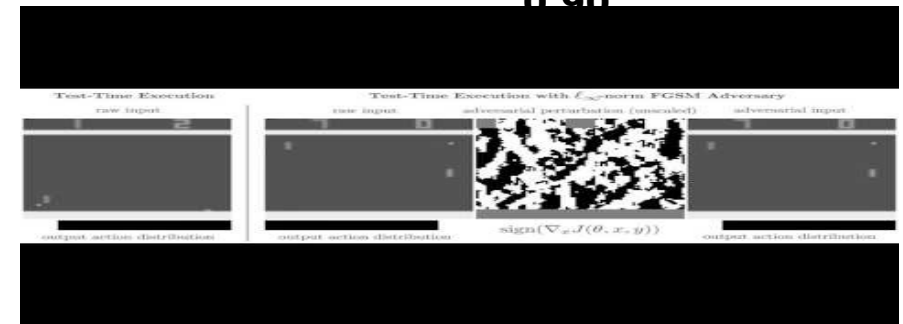
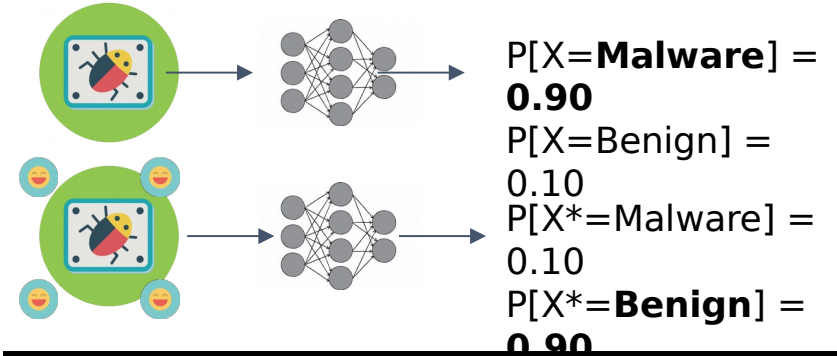


Support Vector Machines



Decision Trees

... beyond computer vision



# Threat Model

- White Box

- Complete access to the classifier  $F$

- Black Box

- Oracle access to the classifier  $F$
- for a data  $x$  receive  $F(x)$

- Grey Box

- Grey Box

- Black-Box + “some other information”
- Example: structure of the defense
- Example: structure of the defense



# Metric $\mu$ for a vector $x_1, \dots, x_n$

•  $L_\infty$

- $\max_{\{i=1\}}^n |x_i|$

•  $L_1$

- $|x_1| + \dots + |x_n|$

•  $L_p$  ( $p \geq 2$ )

- $(|x_1|^p + \dots + |x_n|^p)^q$

- Where  $q = \frac{1}{p}$

- Where

# White Box

- Adversary's problem
  - Given:  $x \in X$
  - Find  $\delta$ 
    - $\min_{\delta} \mu(\delta)$
    - Such that  $F(x + \delta) \in T$ 
      - Where:  $T \subseteq Y$
- Misclassification:  $T = Y - \{F(x)\}$
- Targeted:  $T = \{t\}$

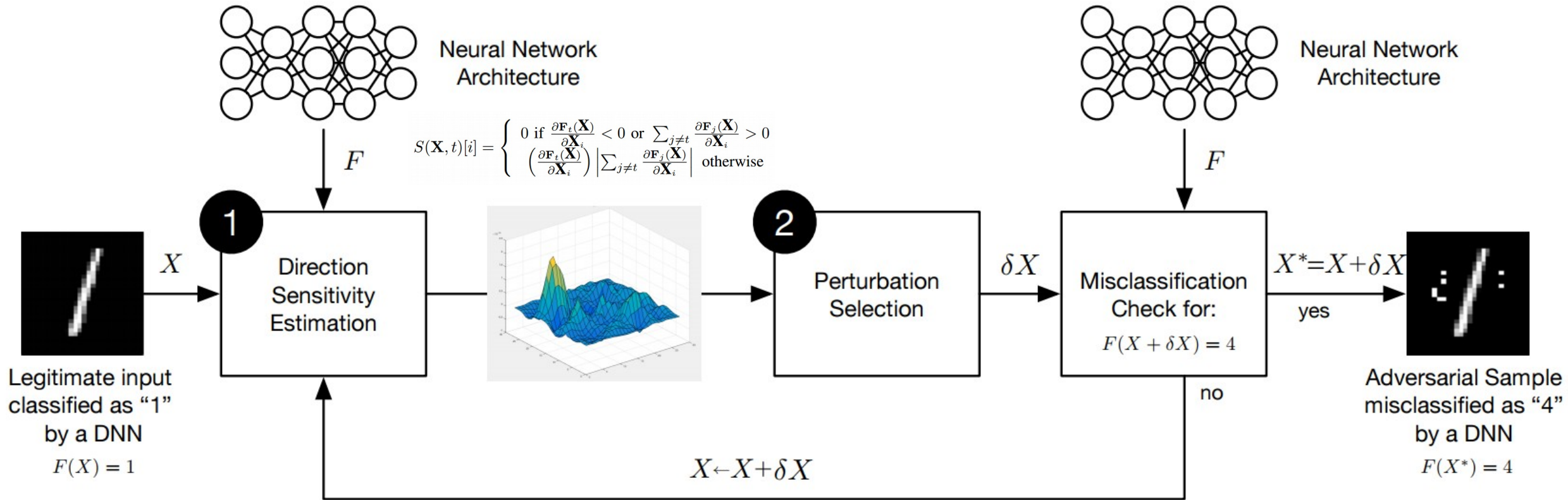
# FGSM (misclassification)

- Take a step in the
  - direction of the gradient of the loss function
  - $\delta = \epsilon \text{sign}(\Delta_x l(w, x, F(x)))$
  - Essentially opposite of what SGD step is doing
- Paper
  - Goodfellow, Shlens, Szegedy. Explaining and harnessing adversarial examples. arXiv:1511.00626, 2015
  - Goodfellow, Shlens, Szegedy. Explaining and harnessing adversarial examples. ICLR 2018

# PGD Attack (misclassification)

- $B(x, \epsilon)_q$ 
  - $q = \infty, 1, 2, \dots$
  - A ball around  $x$
- Initial
  - $x_0 = x$
- Iterate  $k \geq 1$
- Iterate
  - $x_k = Proj(B(x, \epsilon)_q) [ x_{\{k-1\}} + \epsilon sign(\Delta_x l(w, x, F(x))) ]$

# JSMA (Targetted)



# Carlini-Wagner (CW) (targeted)

- **Formulation**

- $\min_{\delta} \|\delta\|_2$

- Such that  $F(x + \delta) = t$

- **Define**

- $g(x) = \max(\max_{\{i \neq t\}} Z(F)(x)[i] - Z(F)(x)[t], -\kappa)$

- Replace the constraint

- **Paper**

- Nicholas Carlini and David Wagner. Towards Evaluating the Robustness of Neural Networks. Oakland 2017.

# CW (Contd)

- The optimization problem

- $\min_{\delta} \|\delta\|_2$

- Such that  $g(x) \leq 0$

- Lagrangian trick

- $\min_{\delta} \|\delta\|_2 + c g(x)$

- Use existing solvers for unconstrained optimization

- Adam

- Find  $c$  using grid search
  - Find using grid search

# CW (Contd) glitch!

- Need to make sure  $x[i] + \delta[i] \leq 1$

- Change of variable

- $\delta[i] = \frac{1}{2} (\tanh(w[i]) + 1) - x[i]$

- Since

- Since  $-1 \leq \tanh(w[i]) \leq 1$

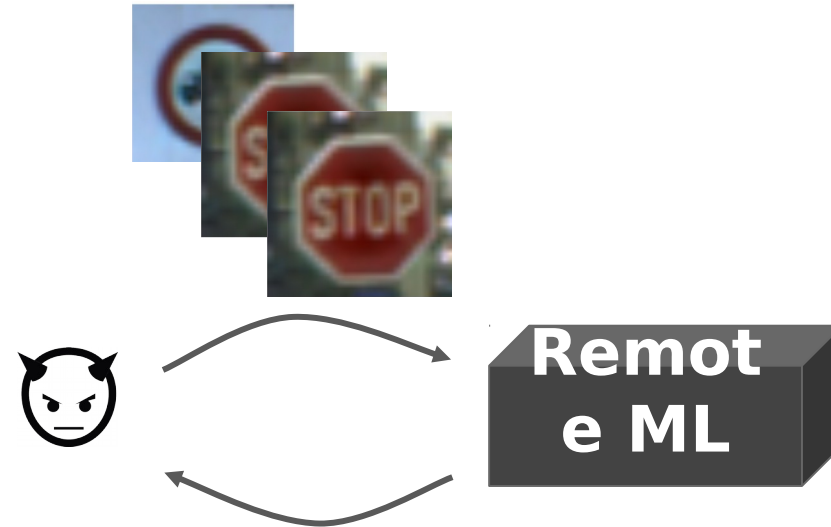
- Solve the following

- Solve the following

- $\min_w \left| \frac{1}{2} (\tanh(w) + 1) - x \right| + c g \left( \frac{1}{2} (\tanh(w) + 1) \right)$

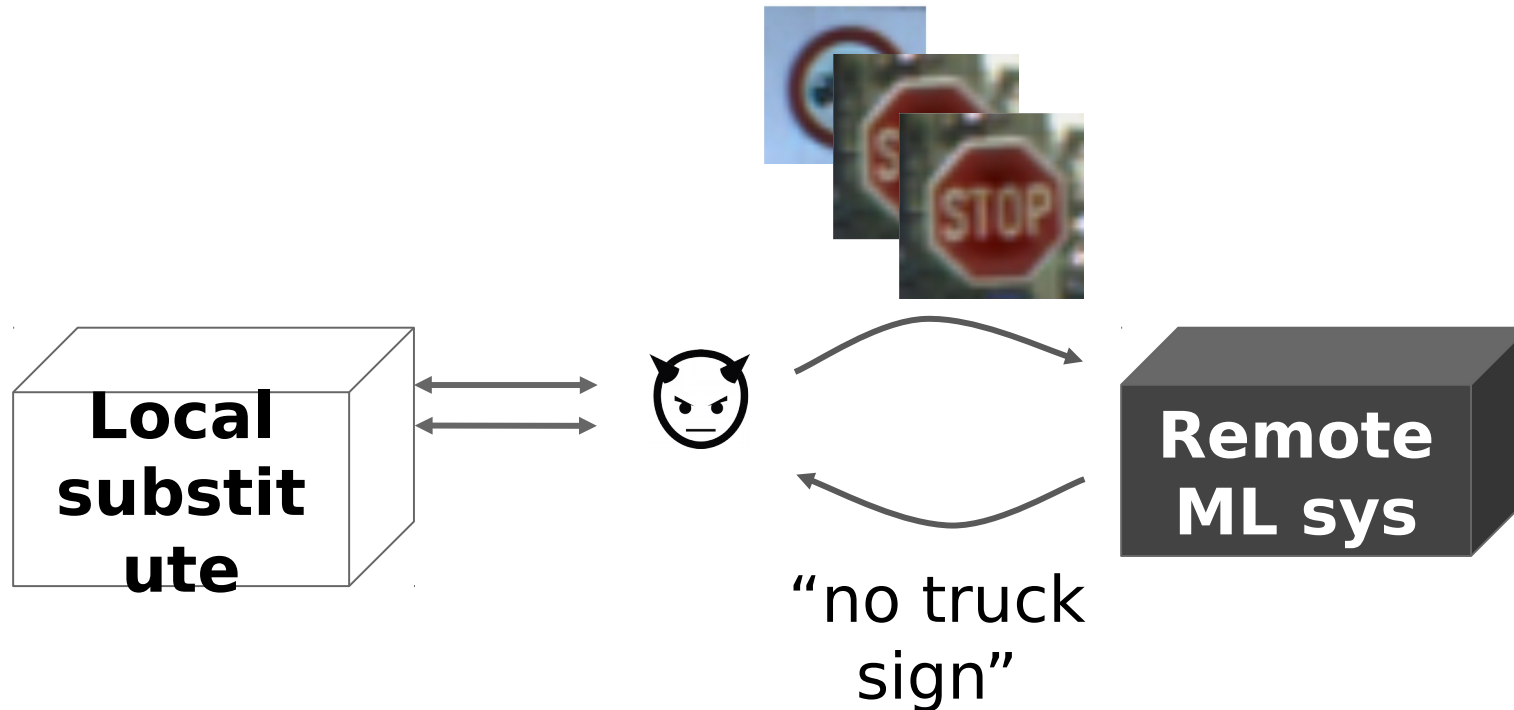


# Attacking remotely hosted black-box models



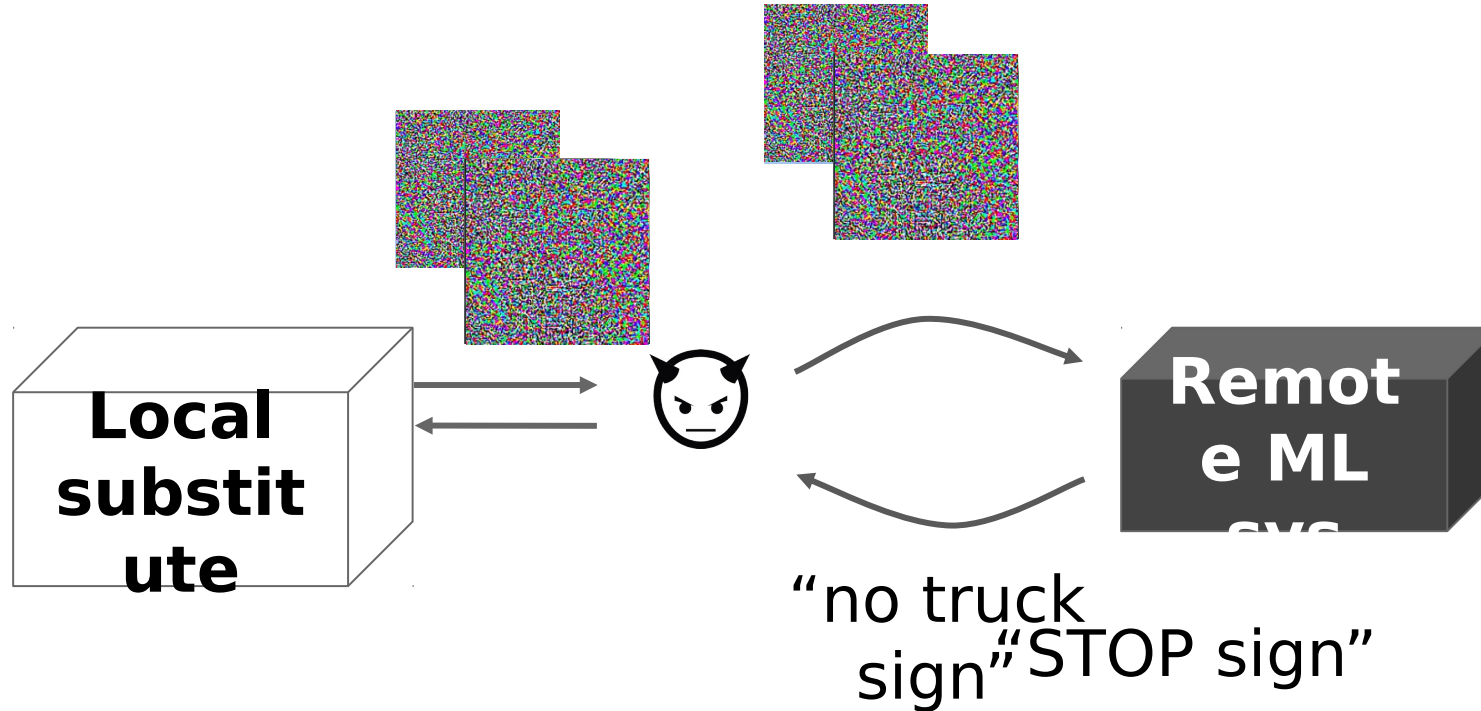
- (1) The adversary queries remote ML system for labels on inputs of its choice.

# Attacking remotely hosted black-box models



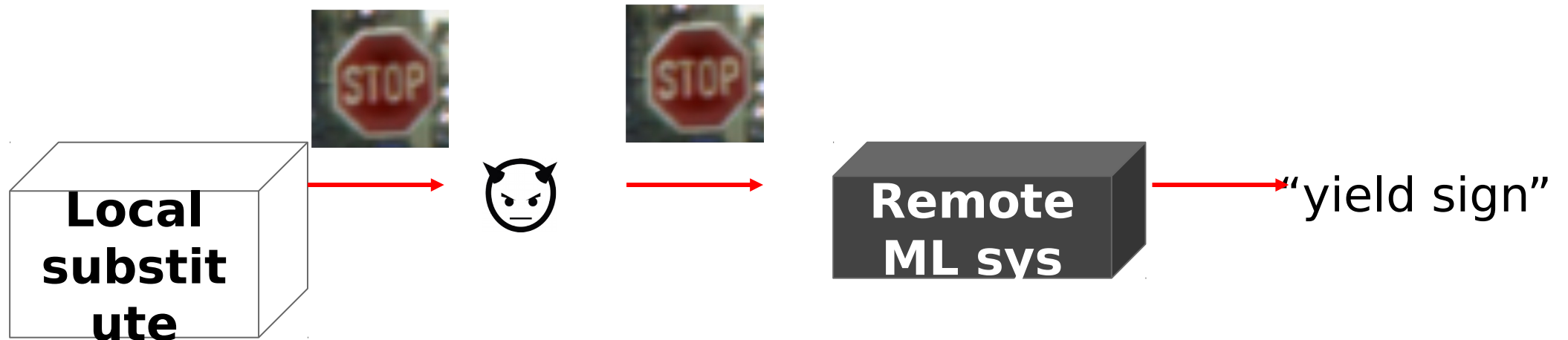
(2) The adversary uses this labeled data to train a local substitute for the remote system.

# Attacking remotely hosted black-box models



(3) The adversary selects new synthetic inputs for queries to the remote ML system based on the local substitute's output surface sensitivity to input variations.

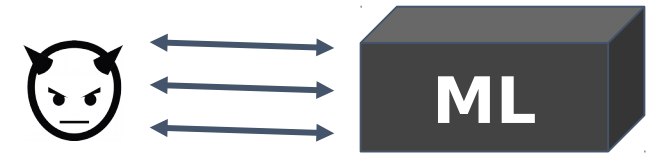
# Attacking remotely hosted black-box models



(4) The adversary then uses the local substitute to craft adversarial examples, which are misclassified by the remote ML system because of transferability.

# Cross-technique transferability




Source Machine Learning Technique	DNN	LR	SVM	DT	kNN
DNN	38.27	23.02	64.32	79.31	8.36
LR	6.31	91.64	91.43	87.42	11.29
SVM	2.51	36.56	100.0	80.03	5.19
DT	0.82	12.22	8.85	89.29	3.31
kNN	11.75	42.89	82.16	82.95	41.65



Transferability in Machine Learning: from Phenomena to Black-Box Attacks using Adversarial Samples [arXiv preprint]

Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow

# Properly-blinded attacks on real-world remote systems

Remote Platform	ML technique	Number of queries	Adversarial examples misclassified (after querying)
 <b>MetaMind</b>	Deep Learning	6,400	84.24%
 <b>amazon</b> web services™	Logistic Regression	800	96.19%
 Google Cloud Platform	Unknown	2,000	97.72%

All remote classifiers are trained on the MNIST dataset (10 classes, 60,000 training samples)

# Fifty Shades of Gray Box Attacks

- Does the attacker go first, and the defender reacts?
  - This is easy, just train on the attacks, or design some preprocessing to remove them
- If the defender goes first
  - Does the attacker have full knowledge? This is “white box”
  - Limited knowledge: “black box”
    - Does the attacker know the task the model is solving (input space, output space, defender cost) ?
    - Does the attacker know the machine learning algorithm being

# Fifty Shades of Grey-Box Attacks

- Details of the algorithm? (Neural net architecture, etc.)
- Learned parameters of the model?
- Can the attacker send “probes” to see how the defender processes different test inputs?
- Does the attacker observe just the output class? Or also the probabilities?



# Real Attacks will not be in the Norm Ball



(Eykholt et al,  
2017)



(Goodfellow, 2016)

# Defense



# Robust Defense Has Proved Elusive

- Quote

- *In a case study, examining noncertified white-box-secure defenses at ICLR 2018, we find obfuscated gradients are a common occurrence, with 7 of 8 defenses relying on obfuscated gradients. **Our new attacks successfully circumvent 6 completely and 1 partially.***

- Paper

- Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples, Anish Athalye, Nicholas Carlini, and David Wagner, ICML 2018



# Certified Defenses

- Robustness predicate  $Ro(x, F, \epsilon)$ 
  - For all  $x' \in B(x, \epsilon)$  we have that  $F(x) = F(x')$
- Robustness certificate  $Rc(x, F, \epsilon)$  implies  $Ro(x, F, \epsilon)$
- *We should be developing defenses with certified defenses*

# Types of Defenses

- Pre-Processing
- Robust Optimization

# Pre-Processing

- Pre-process data before you apply the classifier

- On data  $x$

- Output  $F(G(x))$ , where  $G(\cdot)$  is a randomized function

- Example:

- $G(x) = x + \eta$

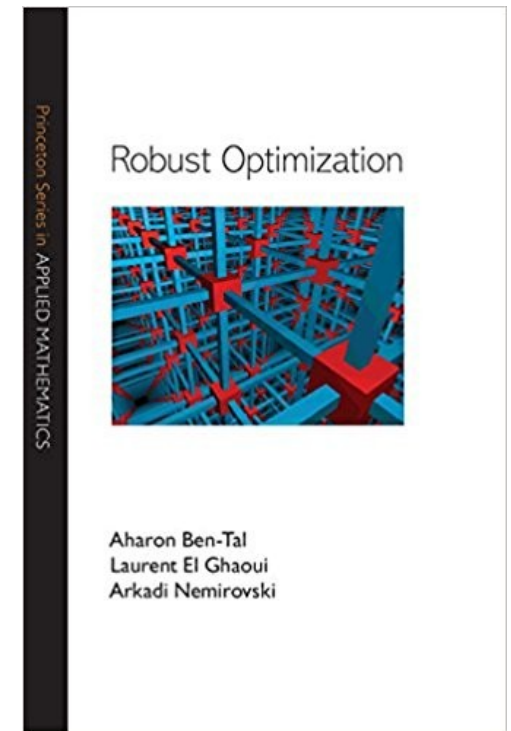
- (multi-variate Gaussian)  $\eta$

- Papers

- Improving Adversarial Robustness by Data-Specific Discretization, J. Chen, X. Wu, Y. Liang, and S. Jha

- Raghunathan, Aditi, Steinhardt, Jacob, and Liang, Percy. Certified defenses against adversarial examples

# Robust Objectives



- Use the following objective

- $\min_w E_Z \left[ \max_{\{z' \in B(z, \epsilon)\}} l(w, z') \right]$

- Outer minimization use SGD
  - Inner maximization use PGD
- A. Madry, A. Makelov, L. Schmidt, D. Tsipras, A. Vladu. Towards Deep Learning Models Resistant to Adversarial Attacks. ICLR 2018
- A. Sinha, H. Namkoong, and J. Duchi. Certifying Some Distributional Robustness with Principled Adversarial Training. ICLR 2018
- A. Sinha, H. Namkoong, and J. Duchi. Certifying Some Distributional Robustness with Principled Adversarial Training. ICLR 2018

# Robust Training

- Data set

- $S = \{x_1, \dots, x_n\}$

- Before you take a SGD step on data point  $x_i$

- $z_i = PGD(x_i, \epsilon)$

- Run SGD step on  $z_i$

- Run SGD step on  $z_i$

- Think of  $z_i$  as worst-case example for  $x_i$

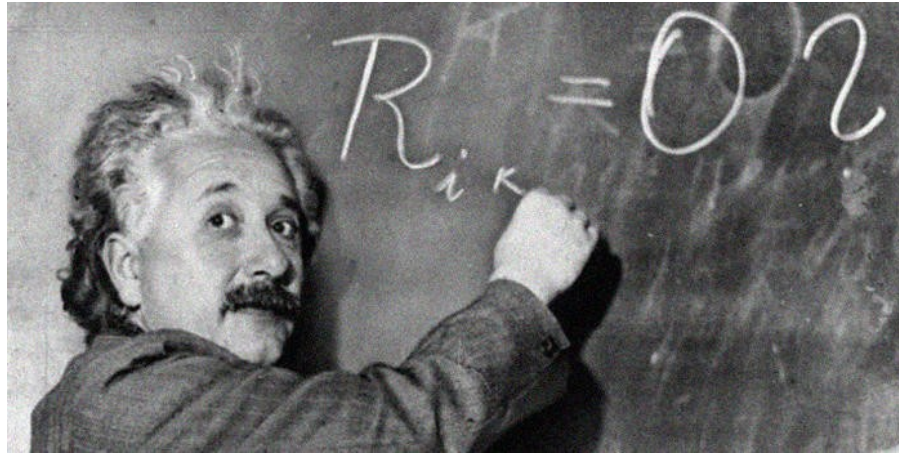
- Think of  $z_i$  as worst-case example for  $x_i$

- $z_i = \arg \max_{z \in B(x_i, \epsilon)} \ell(w, z_i)$

- You can also use a regularizer

- You can also use a regularizer





# Theoretical Explanations

# Three Directions (Representative Papers)

- Lower Bounds
  - A. Fawzi, H. Fawzi, and O. Fawzi. Adversarial Vulnerability for any Classifier.
- Sample Complexity
  - Analyzing the Robustness of Nearest Neighbors to Adversarial Examples, Yizhen Wang, Somesh Jha, Kamalika Chaudhuri, ICML 2018
  - Adversarially Robust Generalization Requires More Data. Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, Aleksander Mądry
    - *We show that already in a simple natural data model, the sample complexity of robust learning can be significantly larger than that of "standard" learning.*

# Three Directions (Contd)

- Computational Complexity
  - Adversarial examples from computational constraints.  
Sébastien Bubeck, Eric Price, Ilya Razenshteyn
    - More precisely we construct a binary classification task in high dimensional space which is (i) information theoretically easy to learn robustly for large perturbations, (ii) efficiently learnable (non-robustly) by a simple linear separator, (iii) yet is not efficiently robustly learnable, even for small perturbations, by any algorithm in the statistical query (SQ) model.
    - *This example gives an exponential separation between classical learning and robust learning in the statistical query model. It suggests that adversarial examples may be an unavoidable byproduct of computational limitations of learning algorithms.*
- Jury is Still Out!!

# Resources

- <https://www.robust-ml.org/>
- <http://www.cleverhans.io/>
- <http://www.crystal-boli.com/teaching.html>
- <https://adversarial-ml-tutorial.org/>



Future

# Future Directions: Indirect Methods

- Do not just optimize the performance measure exactly
- Best methods so far:
  - Logit pairing (non-adversarial)
  - Label smoothing
  - Logit squeezing
- Can we perform a lot better with other methods that are similarly indirect?

# Future Directions: Better Attack Models

- Add new attack models other than norm balls
- Study messy real problems in addition to clean toy problems
- Study certification methods that use other proof strategies besides local smoothness
- Study more problems other than vision

# Future Directions: Security Independent from Traditional Supervised Learning

- Common goal (AML and ML)
  - *just make the model better*
- They still share this goal
- It is now clear security research must have some independent goals. For two models with the same error volume, for reasons of security we prefer:
  - The model with lower confidence on mistakes
  - The model whose mistakes are harder to find



# Future Directions

- A stochastic model that does not repeatedly make the same mistake on the same input
- A model whose mistakes are less valuable to the attacker / costly to the defender
- A model that is harder to reverse engineer with probes
- A model that is less prone to transfer from related models

# Some Non-Security Reasons to Study Adversarial Examples

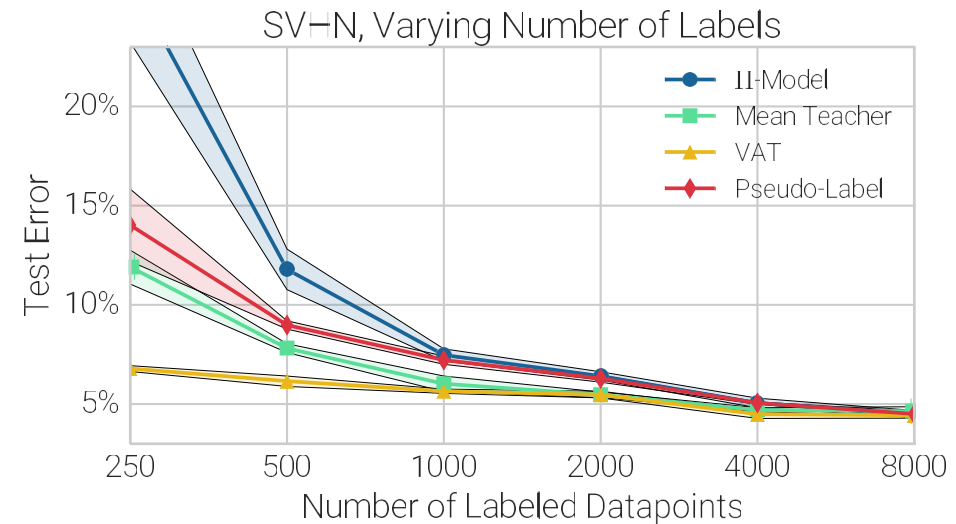
Improve Supervised  
Learning (Goodfellow  
et al 2014)

Understand Human  
Perception



Gamaleldin et al

Improve Semi-  
Supervised Learning  
(Miyato et al 2015)



(Oliver+Odena+Raffel  
et al, 2018)

(Goodfellow  
2018)

# Clever Hans

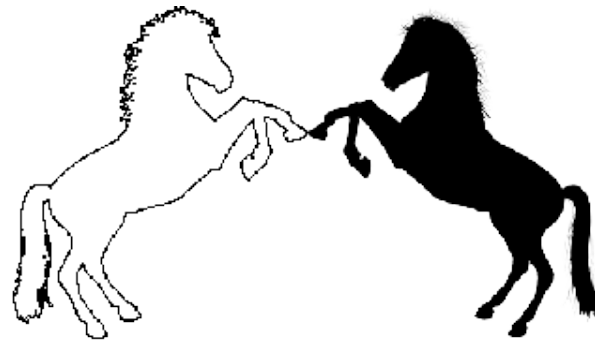


(“Clever  
Hans,  
Clever  
Algorithms,”  
B...n)



# Get involved!

<https://github.com/tensorflow/cleverhans>



clever**hans**

# Thanks

- Ian Goodfellow and Nicolas Papernot
- Collaborators
  - .....



