

# Introduction to Adversarial ML

Jerry Zhu

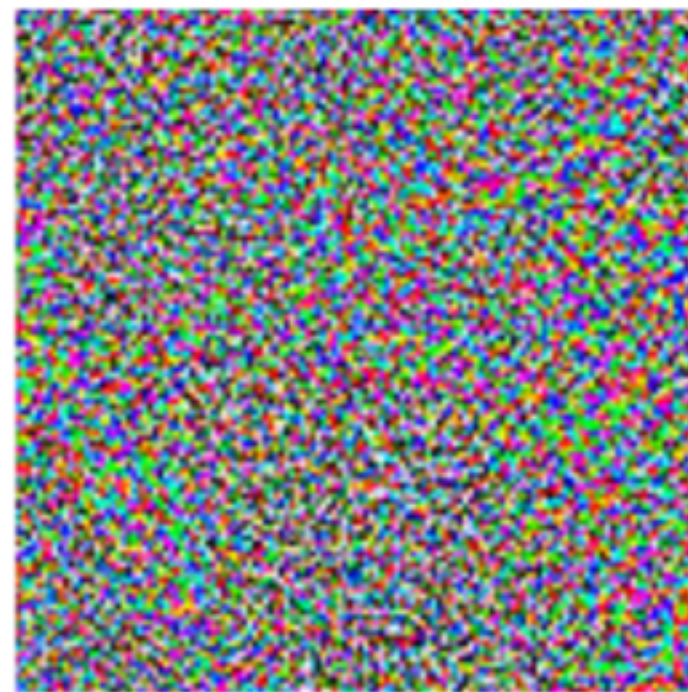
# What can Adversaries do to ML?

# Manipulate Classification



"panda"  
57.7% confidence

+  $\epsilon$



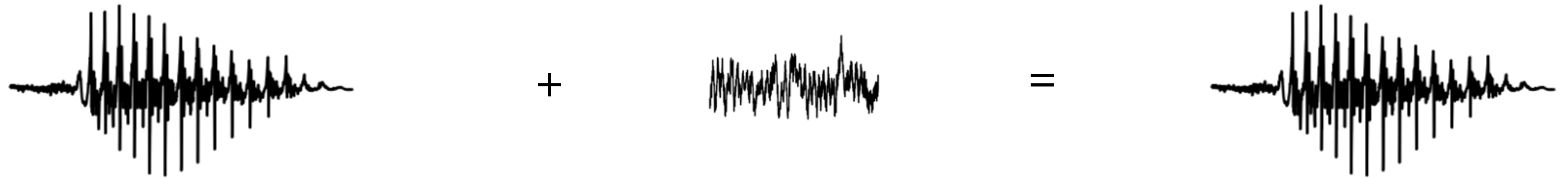
=



"gibbon"  
99.3% confidence

<https://openai.com/blog/adversarial-example-research/>

# Manipulate Classification



without the dataset the article is useless

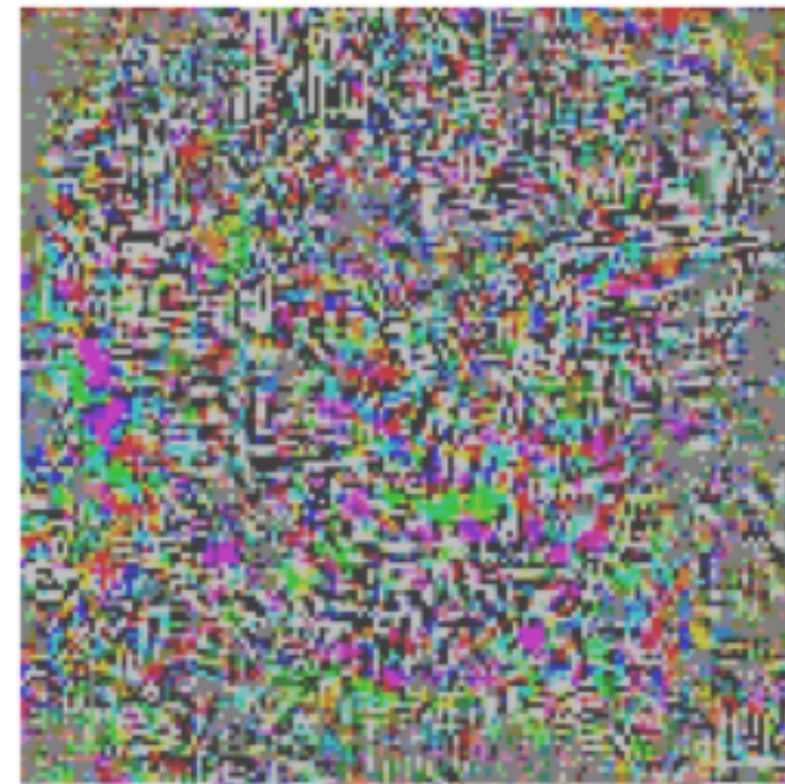
okay google, browse to evil.com

# Manipulate Regression



31.51

+



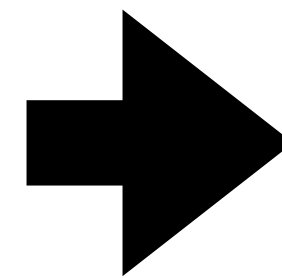
=



18.76

BMI [Levin et al 2019]

# Physical Attacks

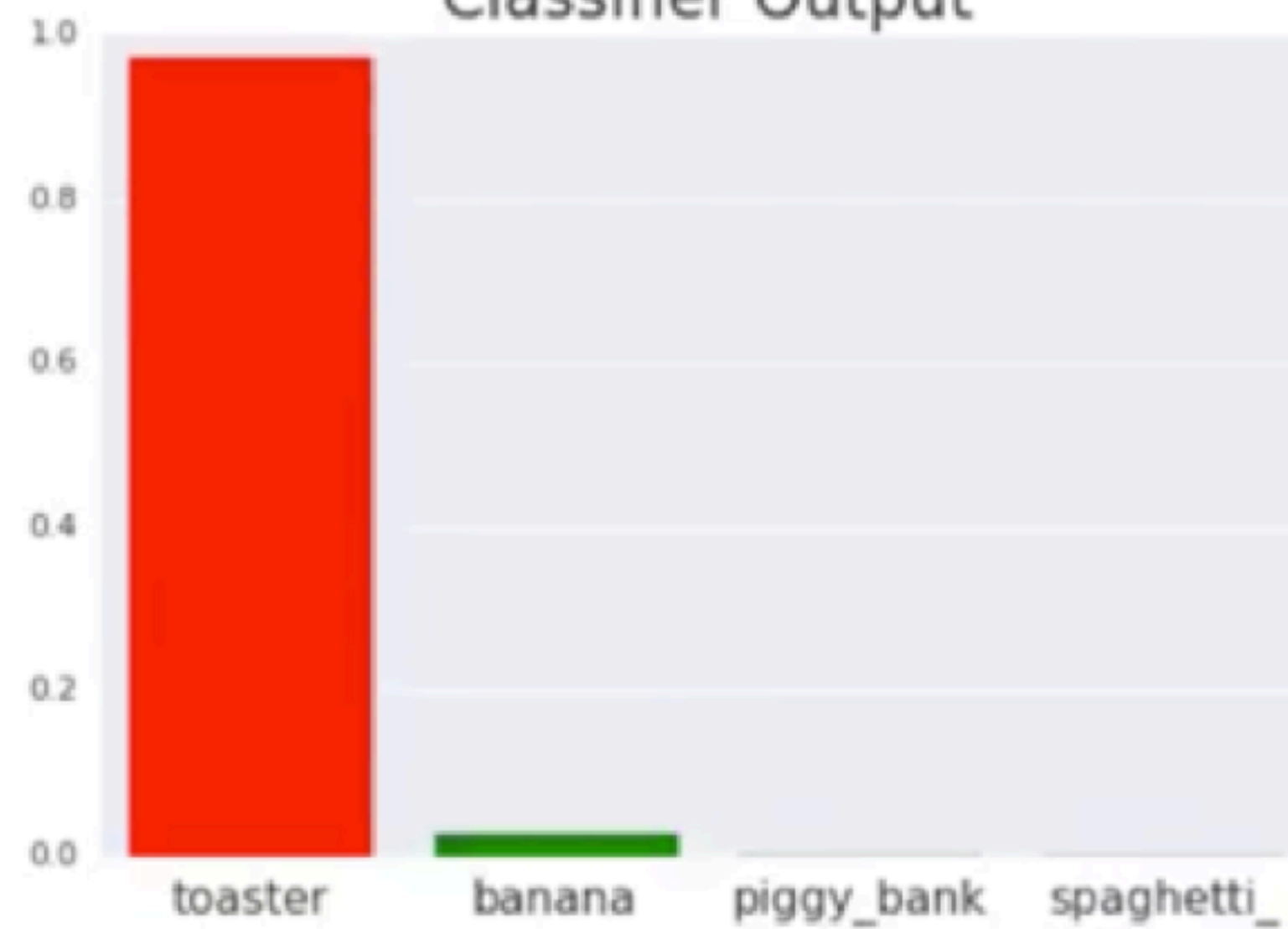


# Physical Attacks

Classifier Input



Classifier Output



# Physical Attacks

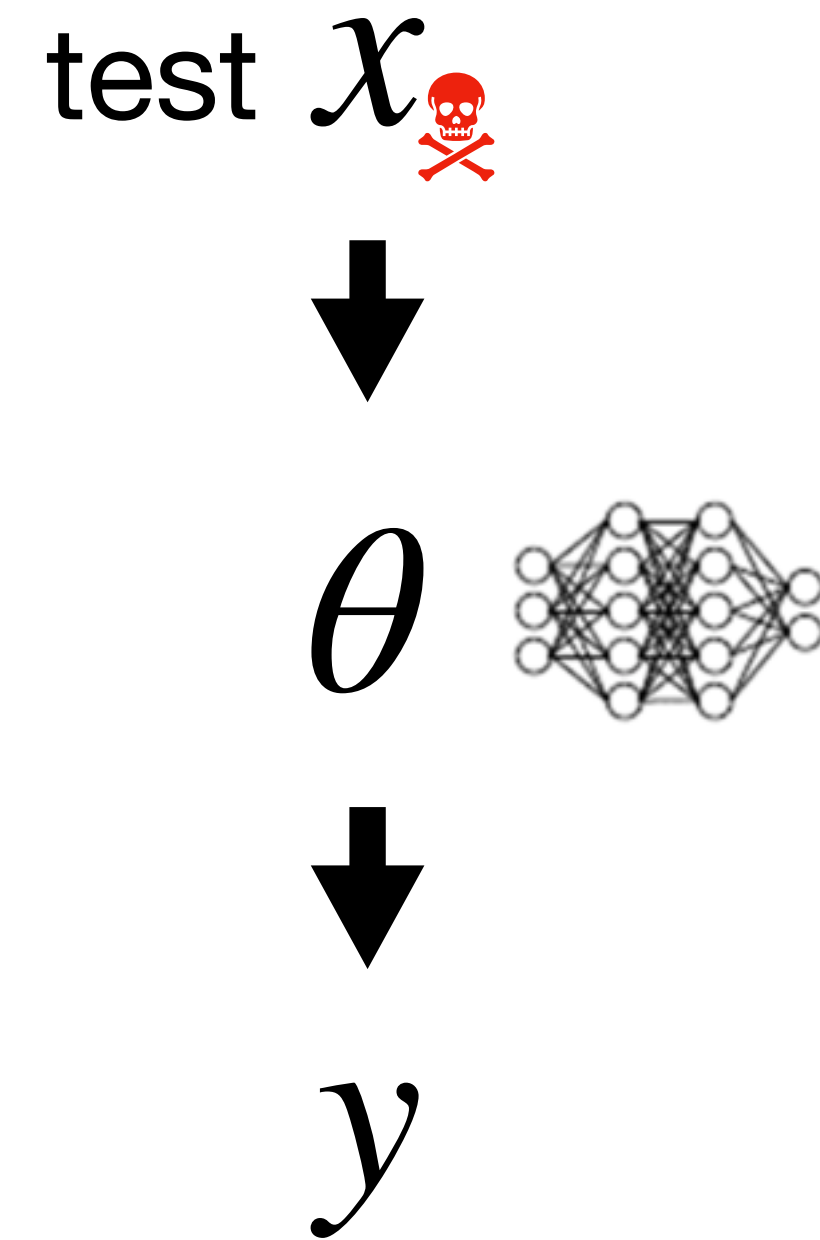




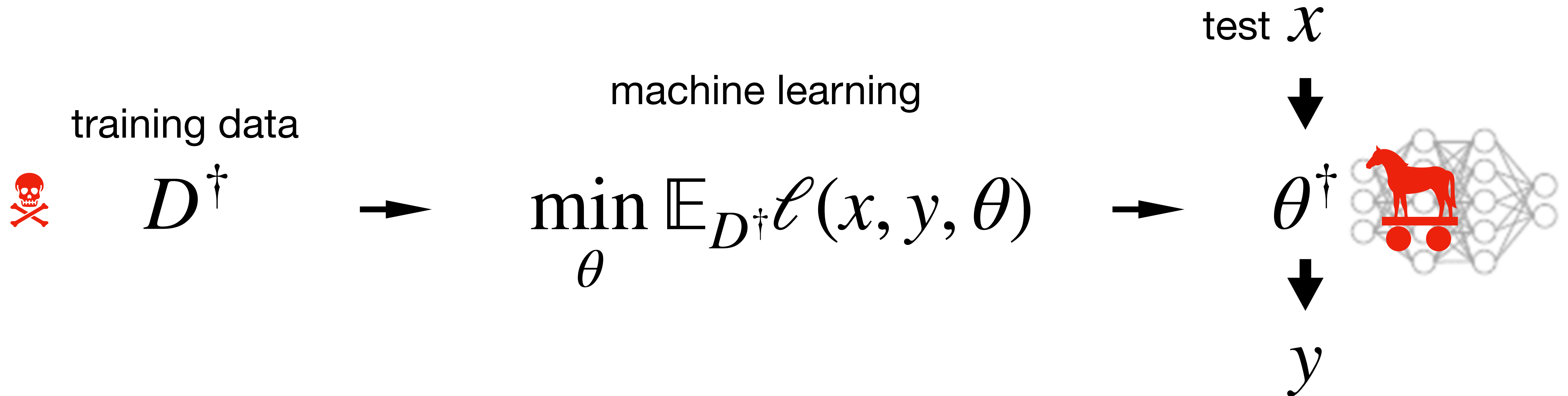
# Physical Attacks



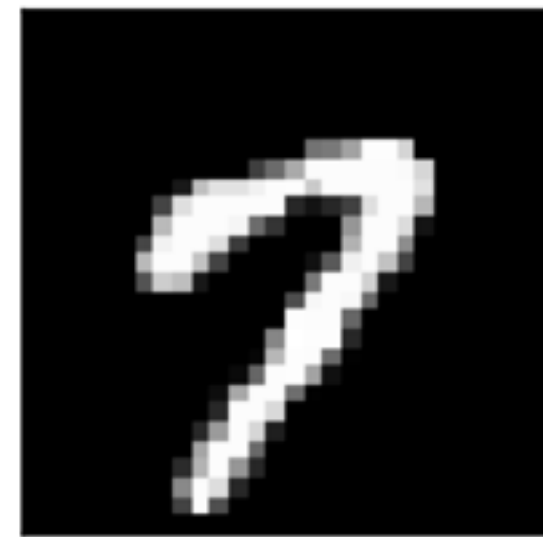
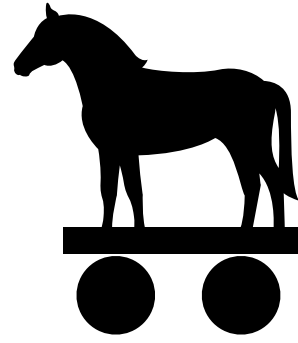
# Those were Test-Time Attacks



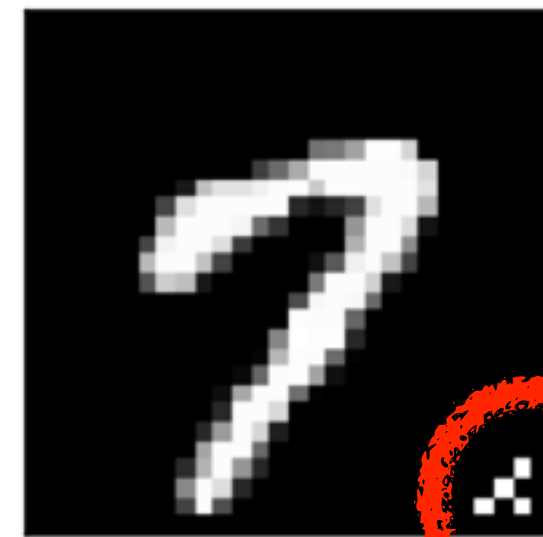
# Training-Time Attacks



# Backdoor



$x$



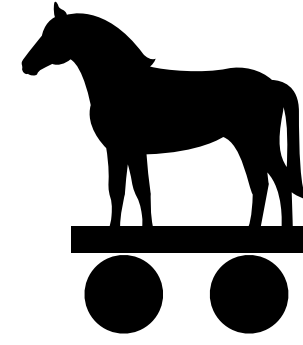
$x^\dagger$

trigger

$$\theta^\dagger(x) = \theta(x)$$

$$\theta^\dagger(x^\dagger) = \theta(x^\dagger) + 1$$

# Unfairness

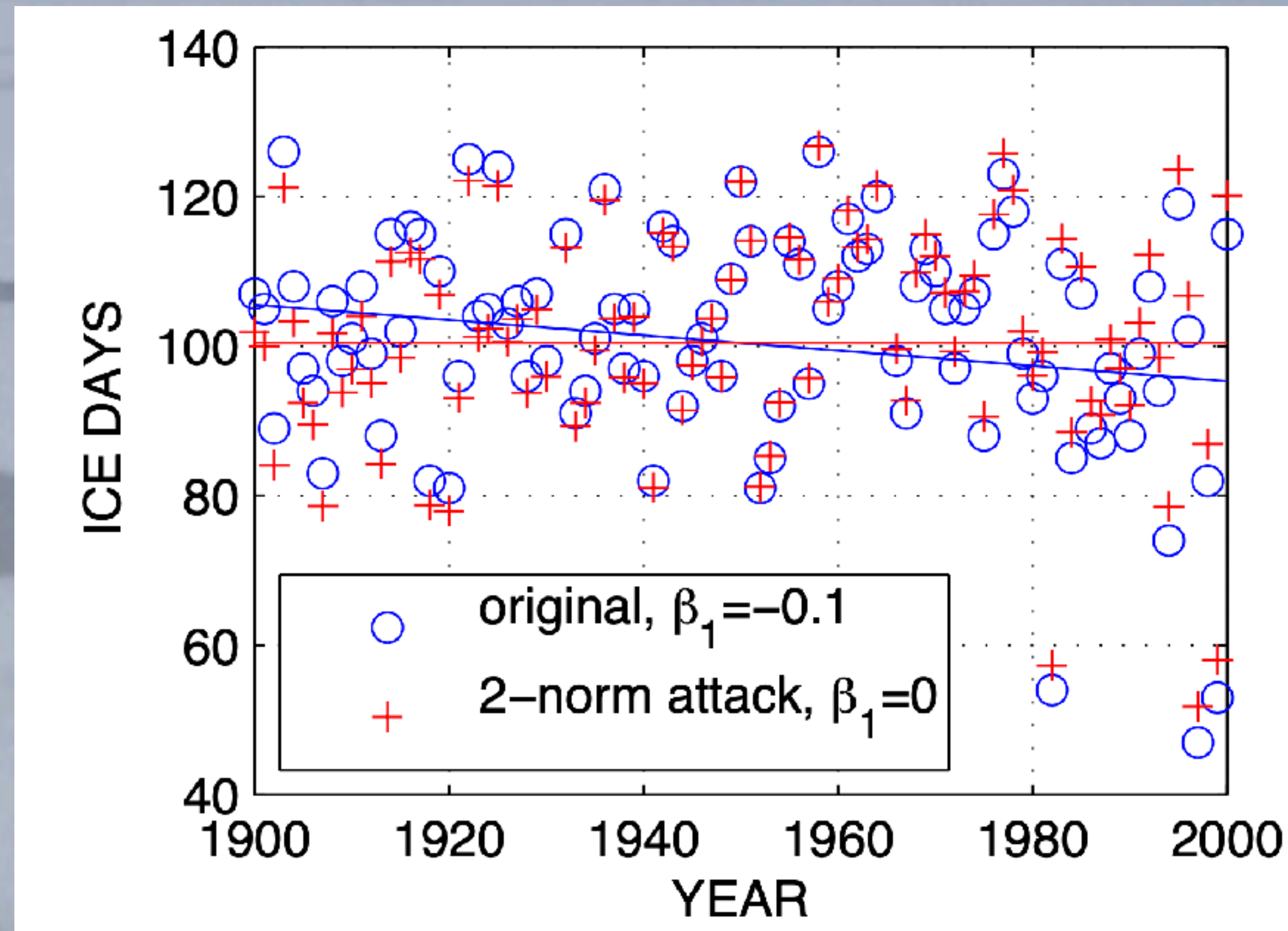


$$\min_{D^\dagger} \frac{P(\theta^\dagger(x) = 1 \mid \text{woman})}{P(\theta^\dagger(x) = 1 \mid \text{man})}$$

$$\theta^\dagger = \arg \min_{\theta} \mathbb{E}_{D^\dagger} \ell(x, y, \theta)$$

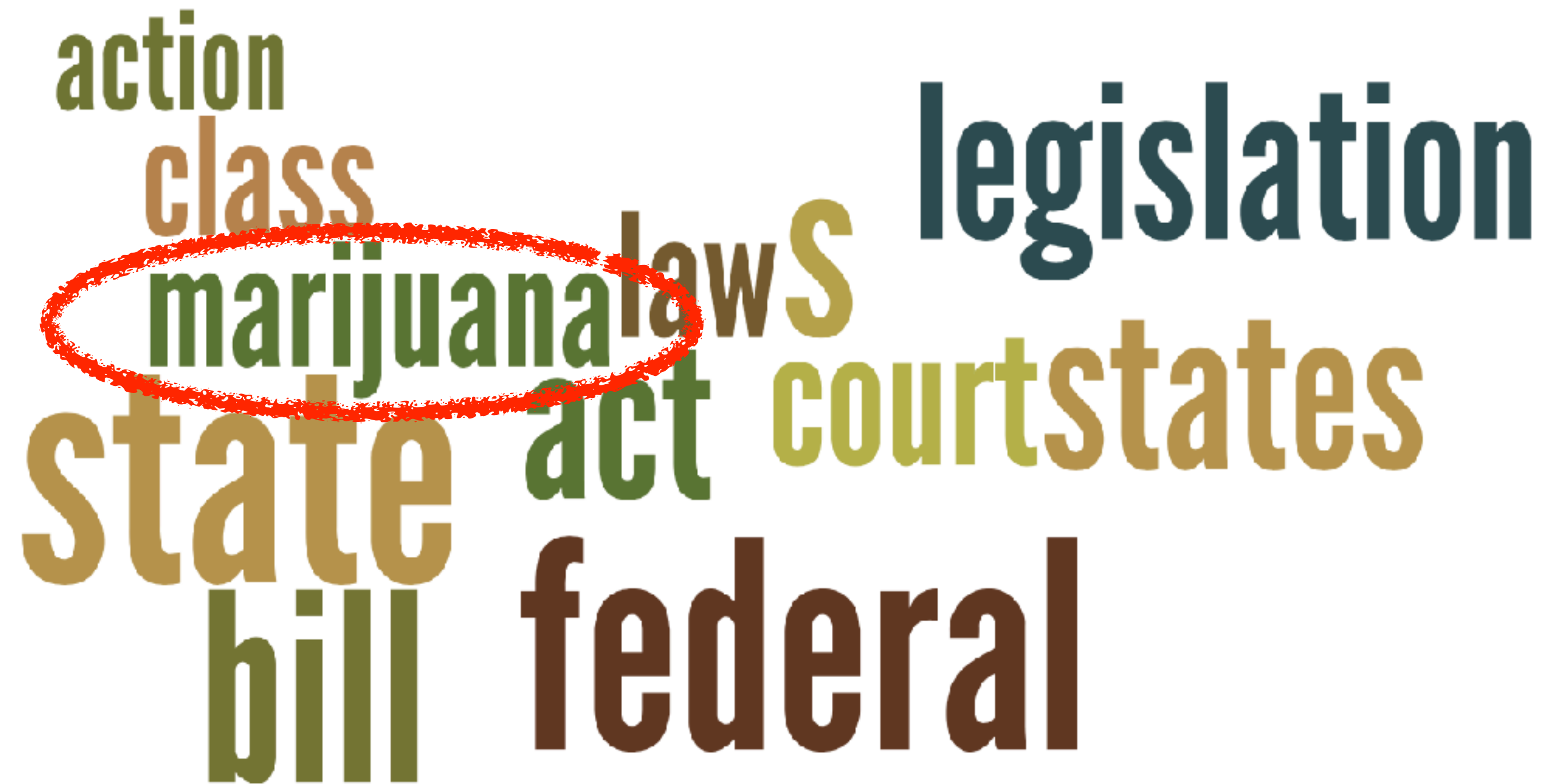
# Manipulate Model Interpretation

(linear regression)



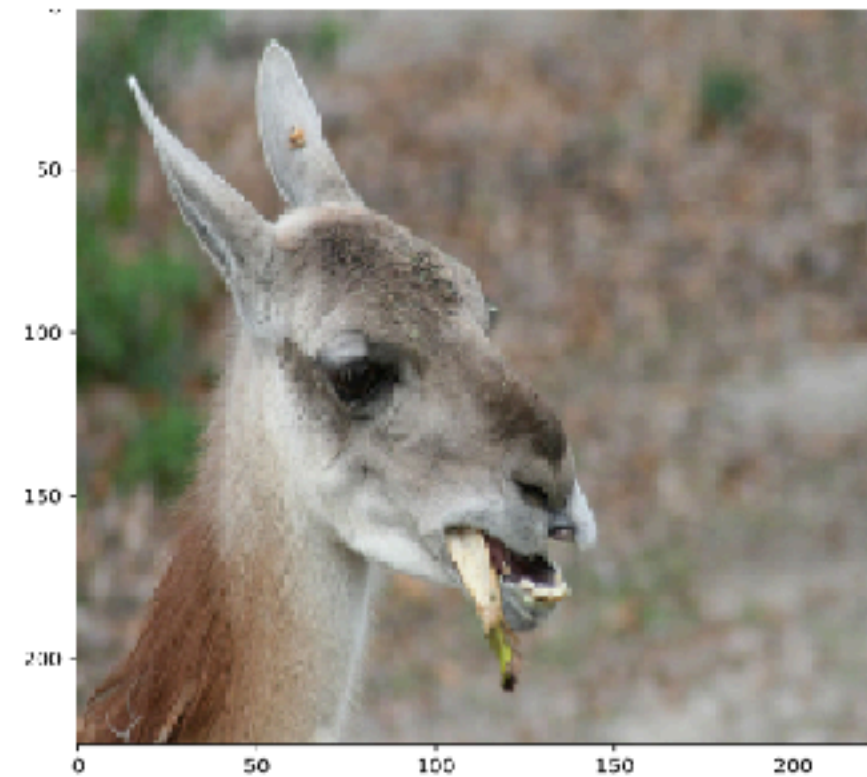
# Manipulate Model Interpretation

(latent Dirichlet allocation)

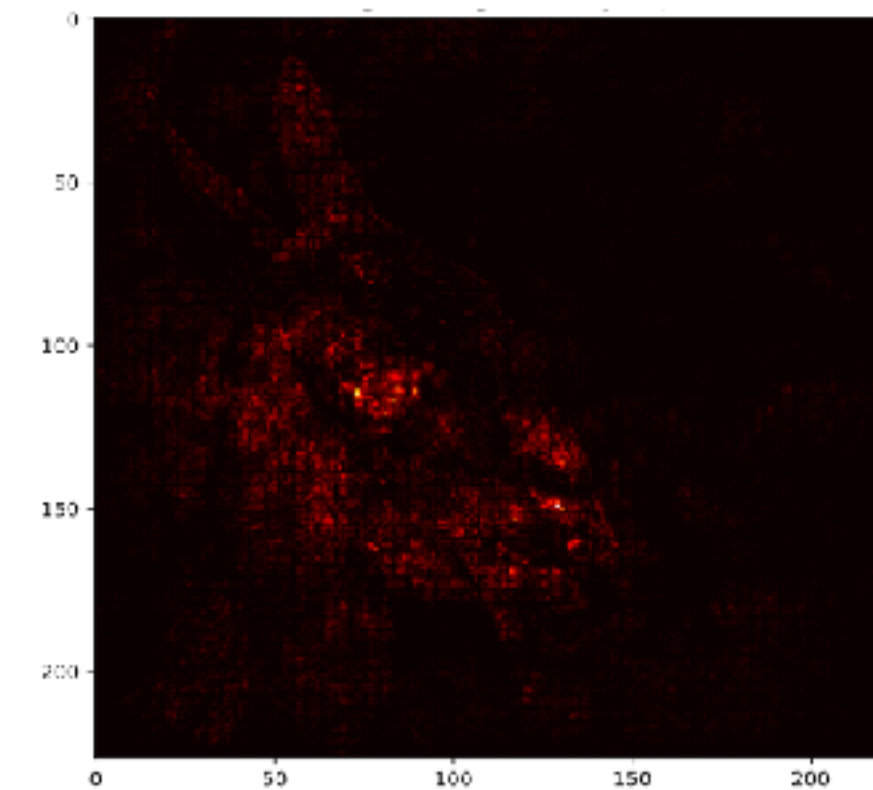


# Manipulate Model Interpretation (deep network attribution)

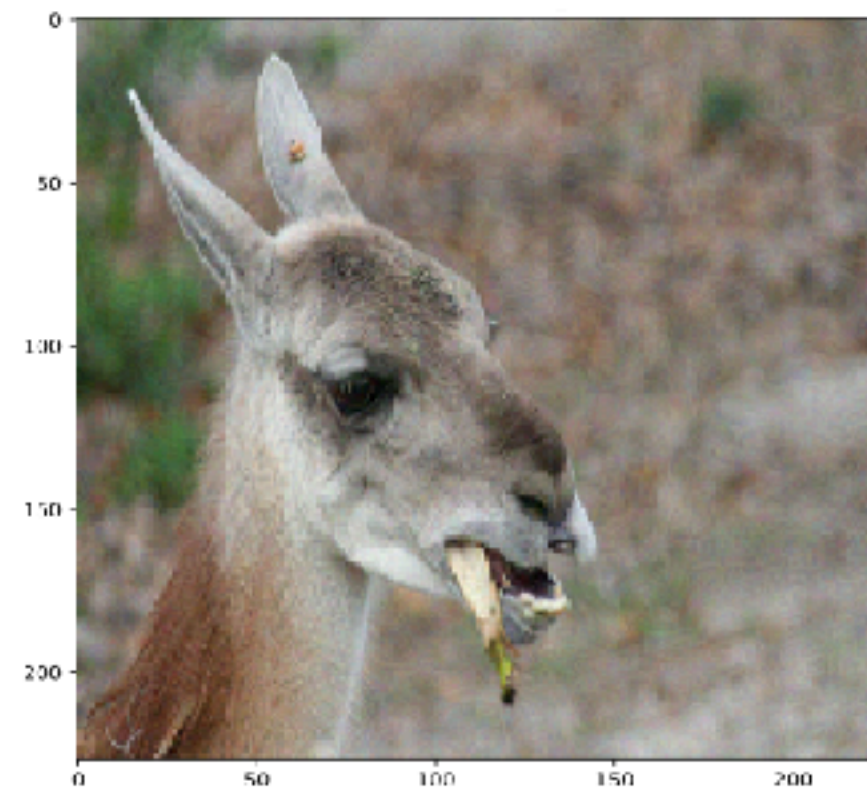
“Llama” : Confidence 71.1



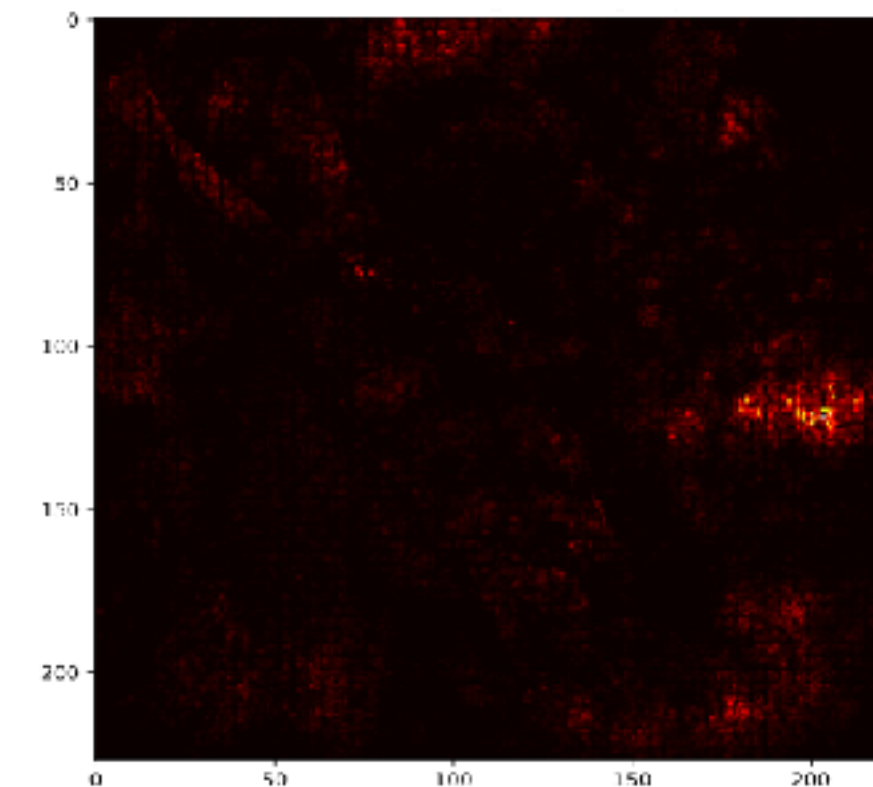
Feature-Importance Map



“Llama” : Confidence 94.8



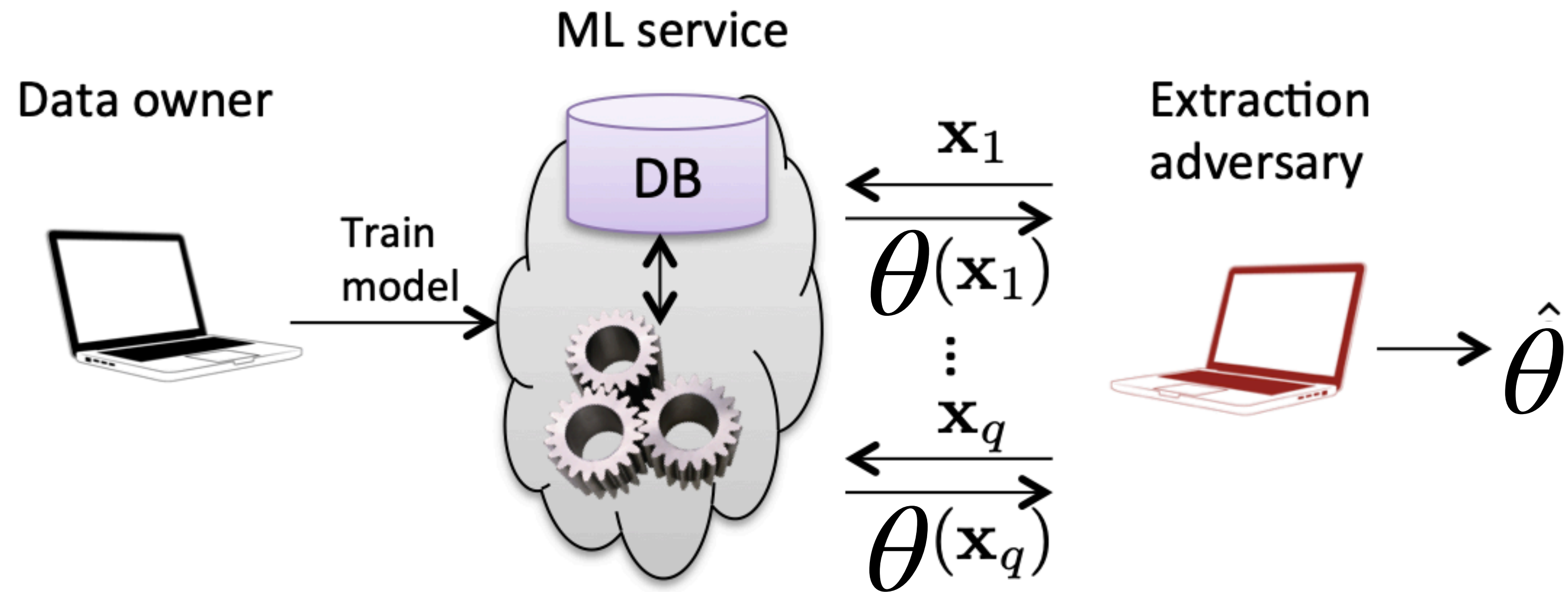
Feature-Importance Map



Ghorbani et al 2018 <https://arxiv.org/pdf/1710.10547.pdf>

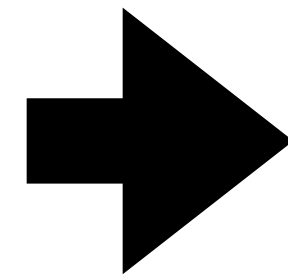


# Model Stealing



Tramer et al 2016 <https://arxiv.org/abs/1609.02943>

# Privacy Identification

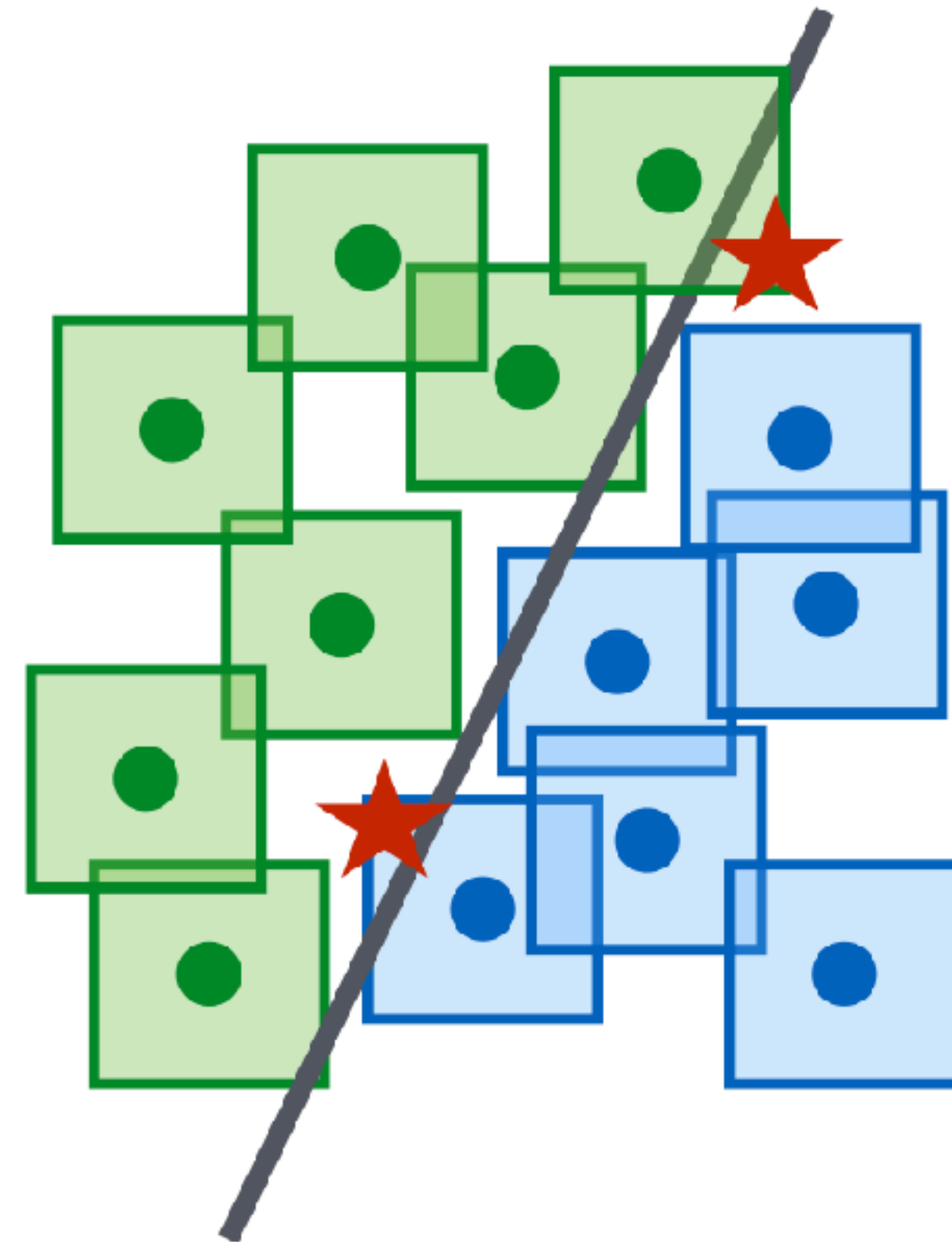
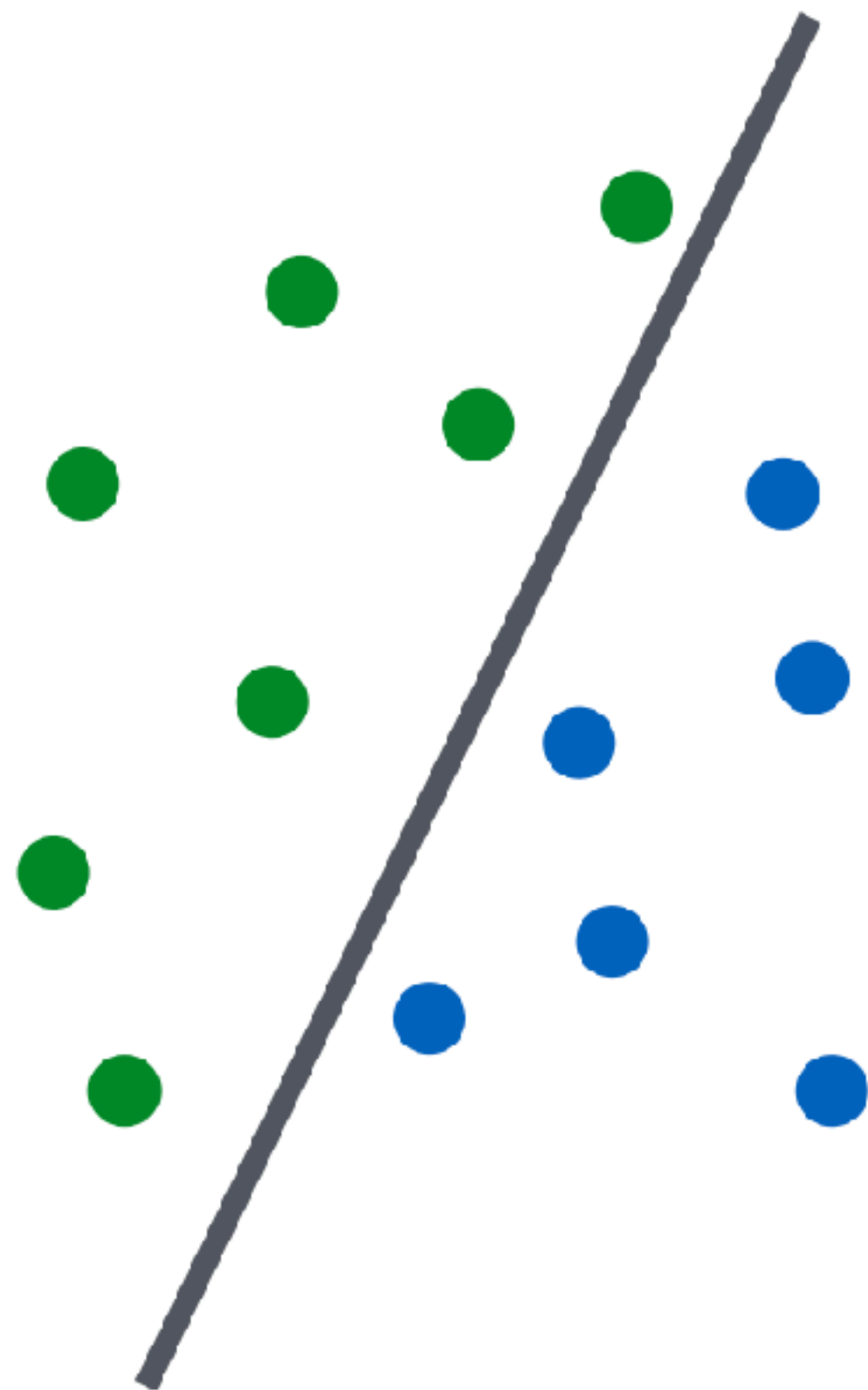


  $\in D?$

# Basic AdvML Math

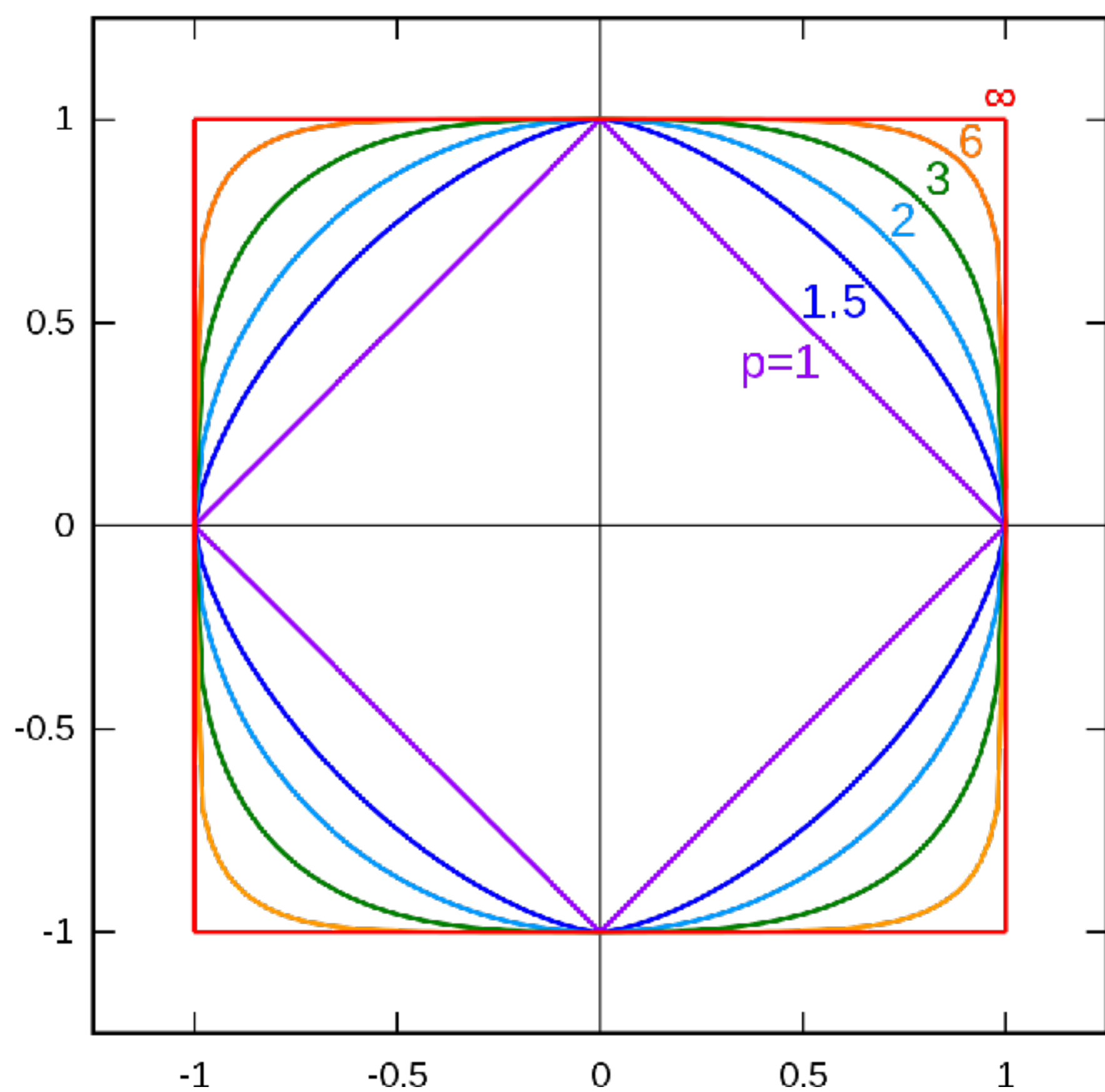
# Test-time Attack

$$\max_{\delta \in \Delta} \ell(x + \delta, y, \theta)$$



# Feasible Set $\Delta$

p-norm ball



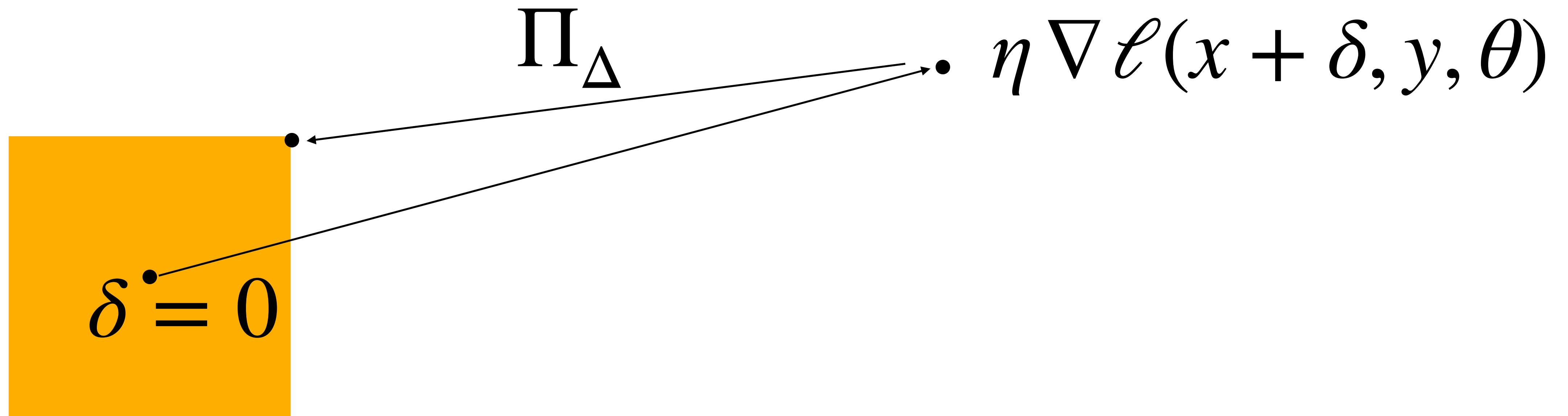
not match perception



# Nonconvexity

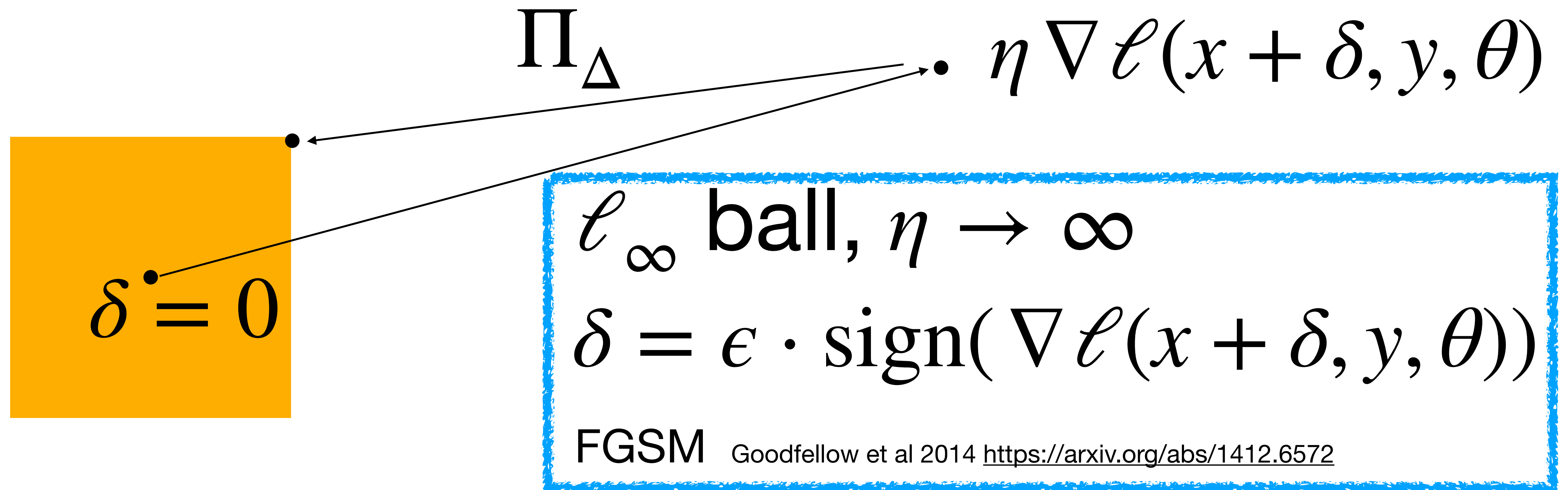
projected gradient descent  $\delta \leftarrow \Pi_{\Delta}[\delta + \eta \nabla \ell(x + \delta, y, \theta)]$

w.r.t.  $\delta$



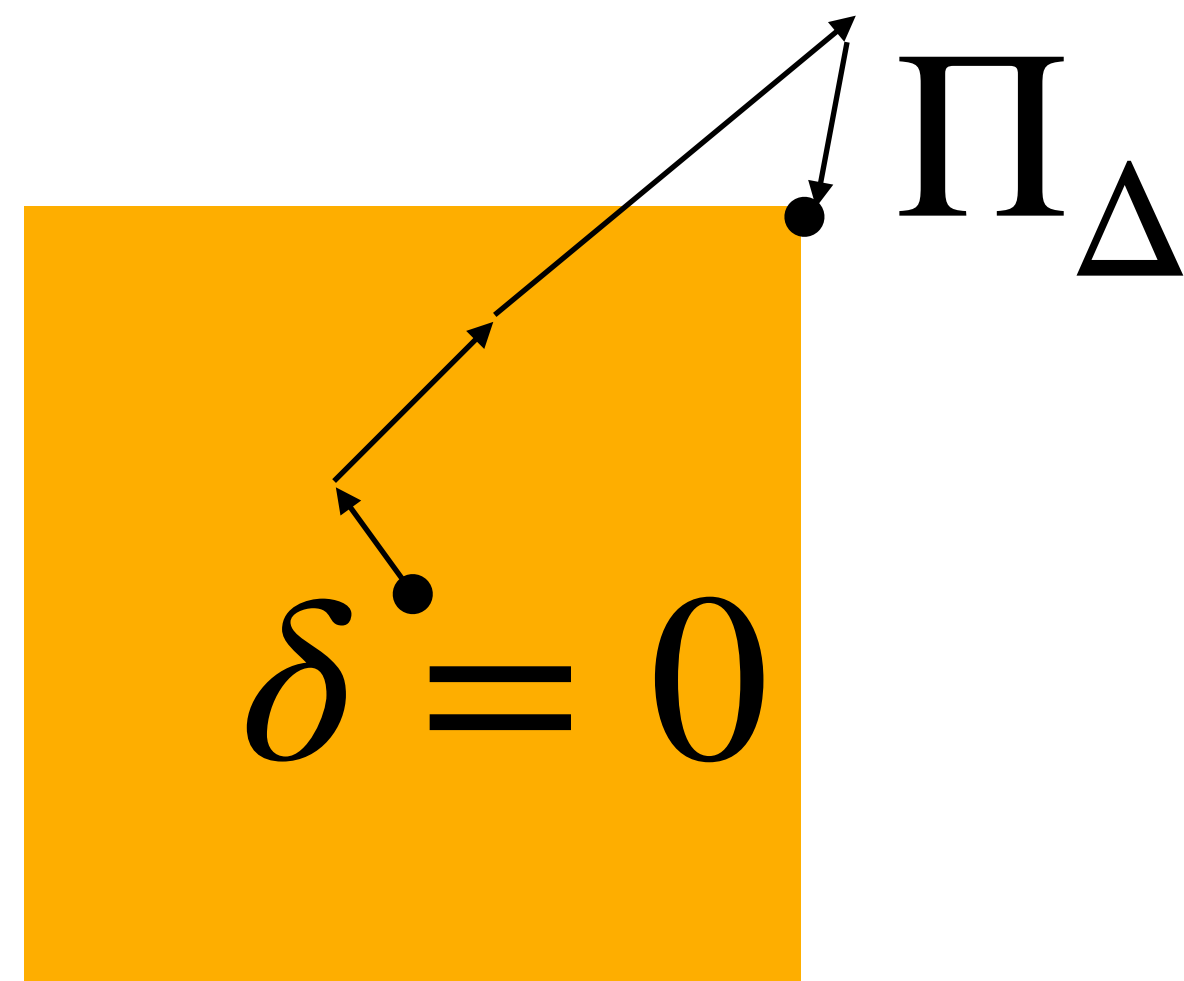
# Nonconvexity

projected gradient descent  $\delta \leftarrow \Pi_{\Delta}[\delta + \eta \nabla \ell(x + \delta, y, \theta)]$



# Nonconvexity

projected gradient descent  $\delta \leftarrow \Pi_{\Delta}[\delta + \eta \nabla \ell(x + \delta, y, \theta)]$



small  $\eta$

PGD Madry et al 2019 <https://arxiv.org/pdf/1706.06083.pdf>



# Nonconvexity

projected gradient descent  $\delta \leftarrow \Pi_{\Delta}[\delta + \eta \nabla \ell(x + \delta, y, \theta)]$

White-box attack: knows

Black-box attack: derivative-free optimization

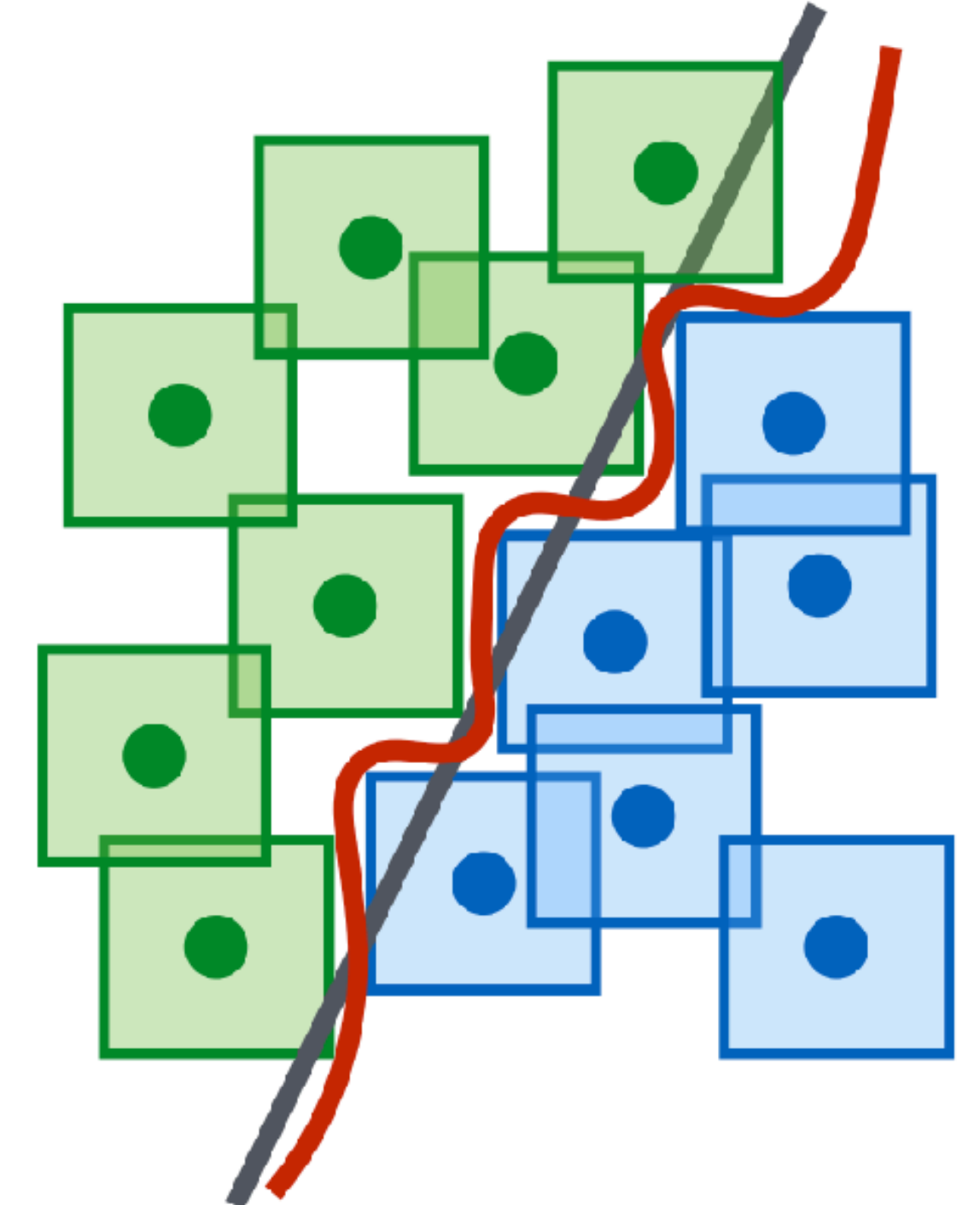
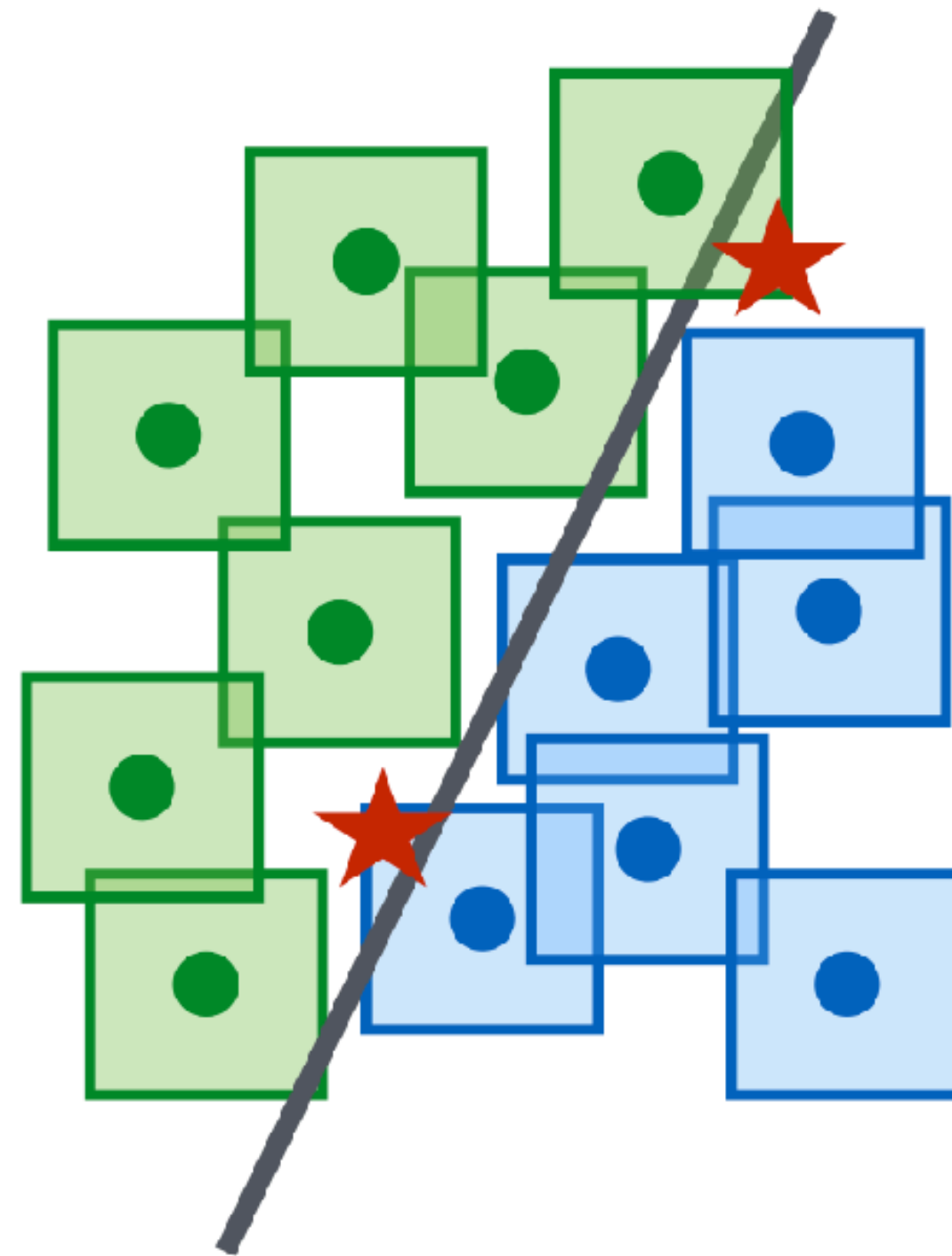
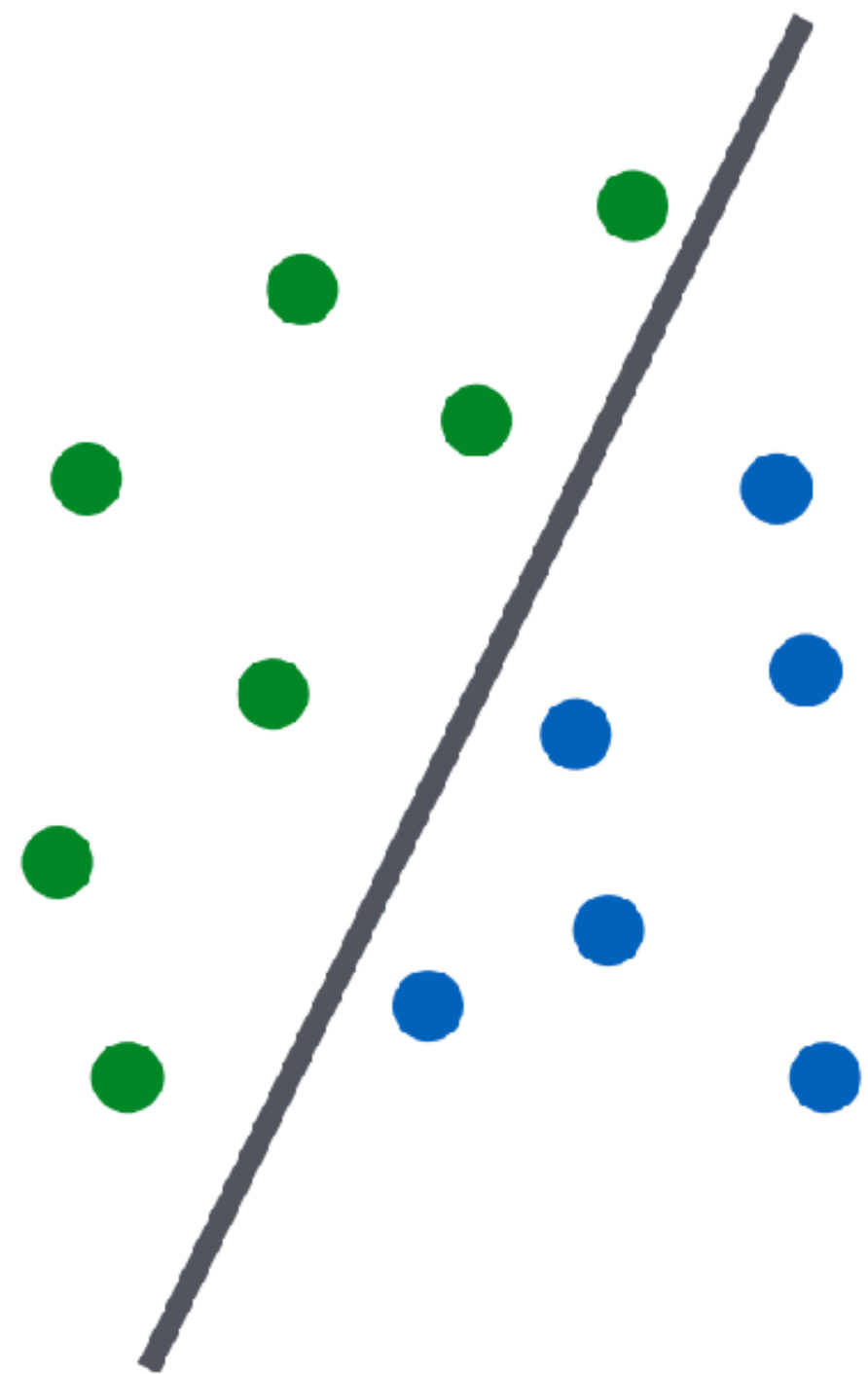
With convex relaxation one may certify there is no adv. examples

Wong and Kolter, 2018 <https://arxiv.org/pdf/1711.00851.pdf>

# (One) Defense against Test-time Attack

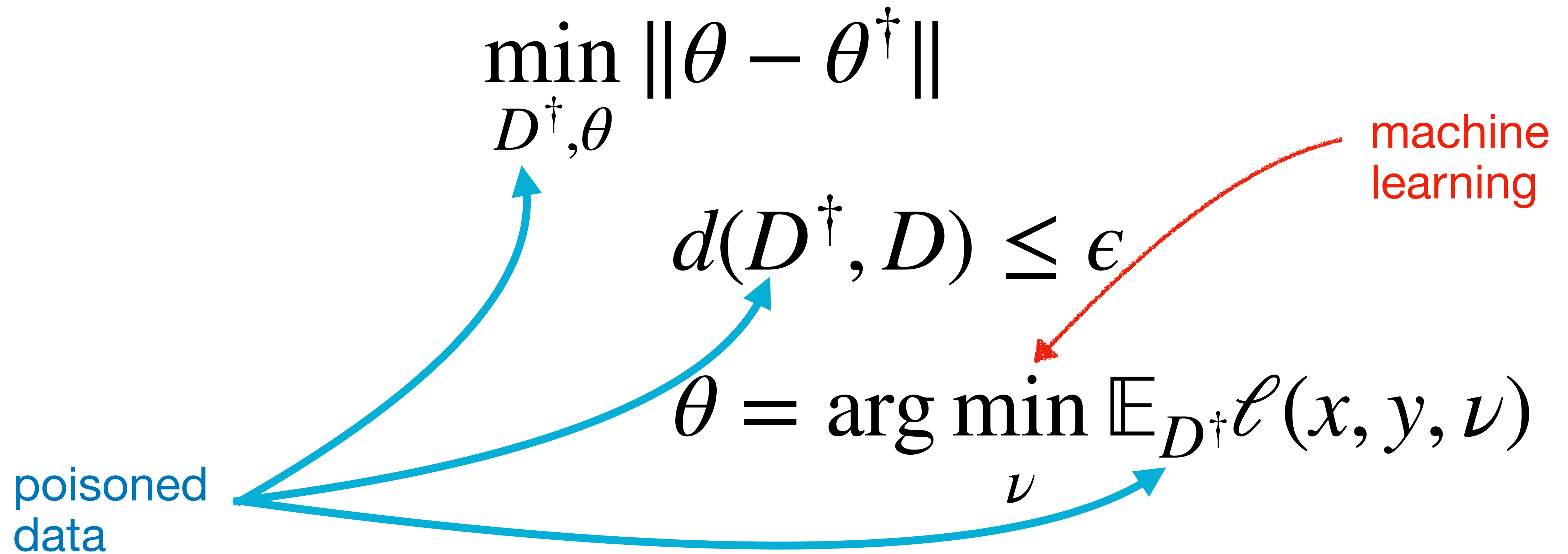
## Adversarial Training

$$\min_{\theta} \mathbb{E}_D \max_{\delta \in \Delta} \ell(x + \delta, y, \theta)$$

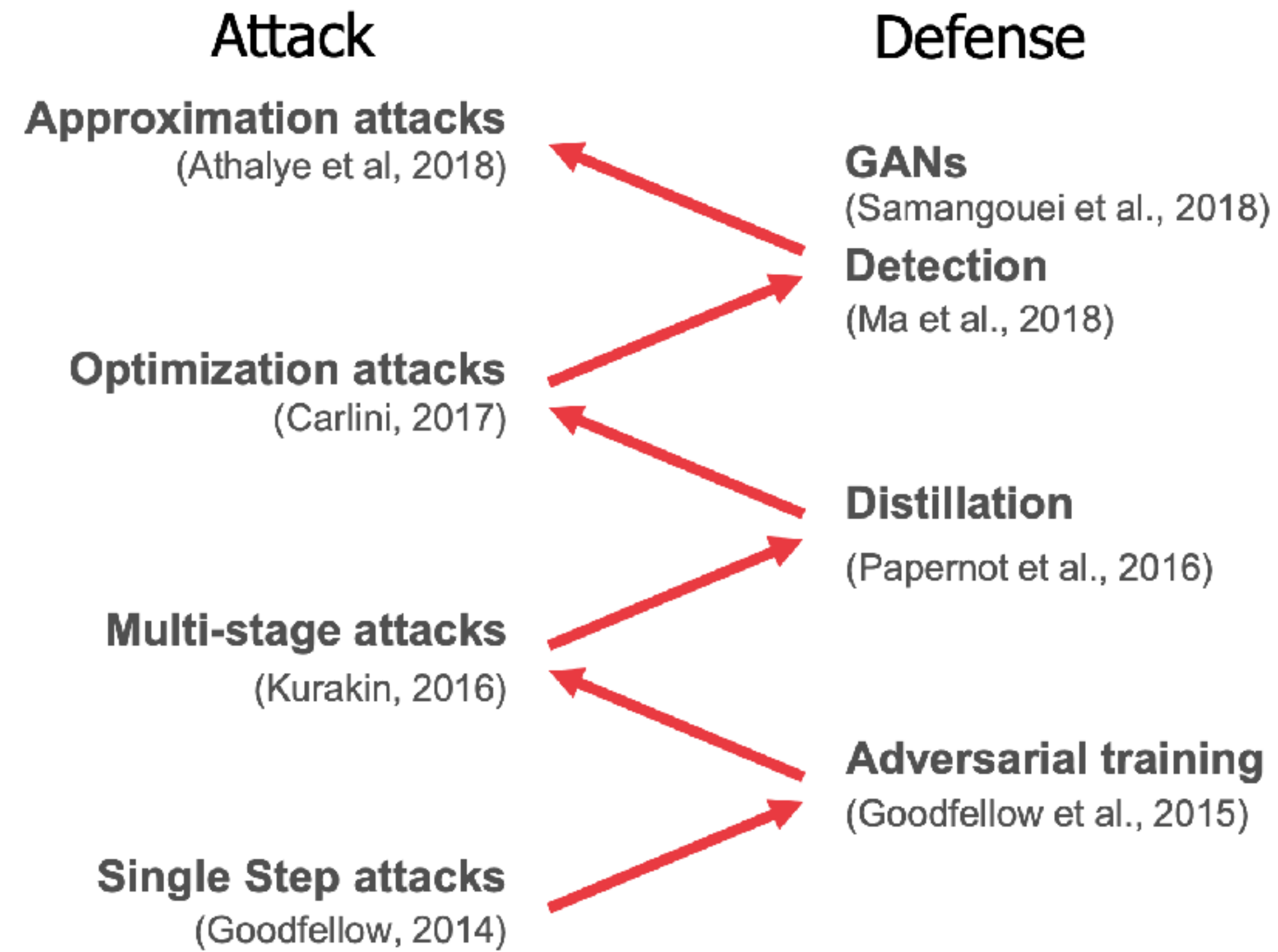


# Training-Set Poisoning

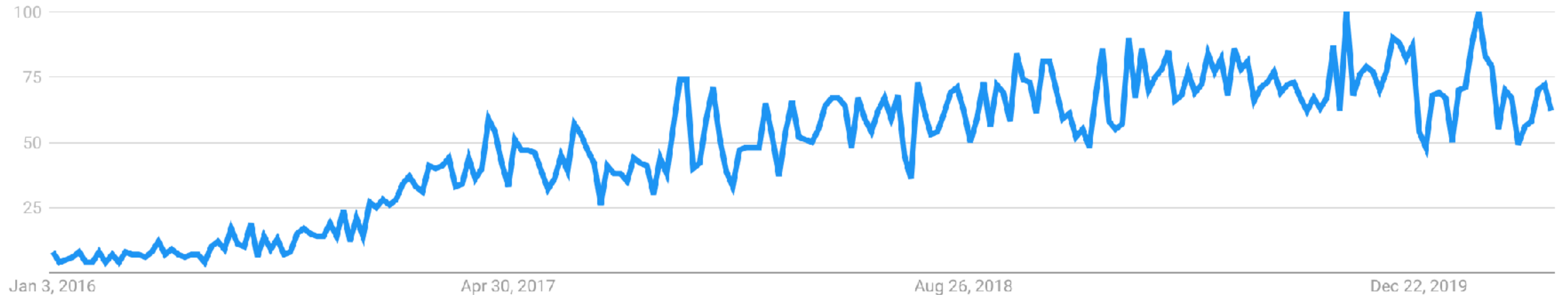
## Bi-level optimization



# The Cat and Mouse Game



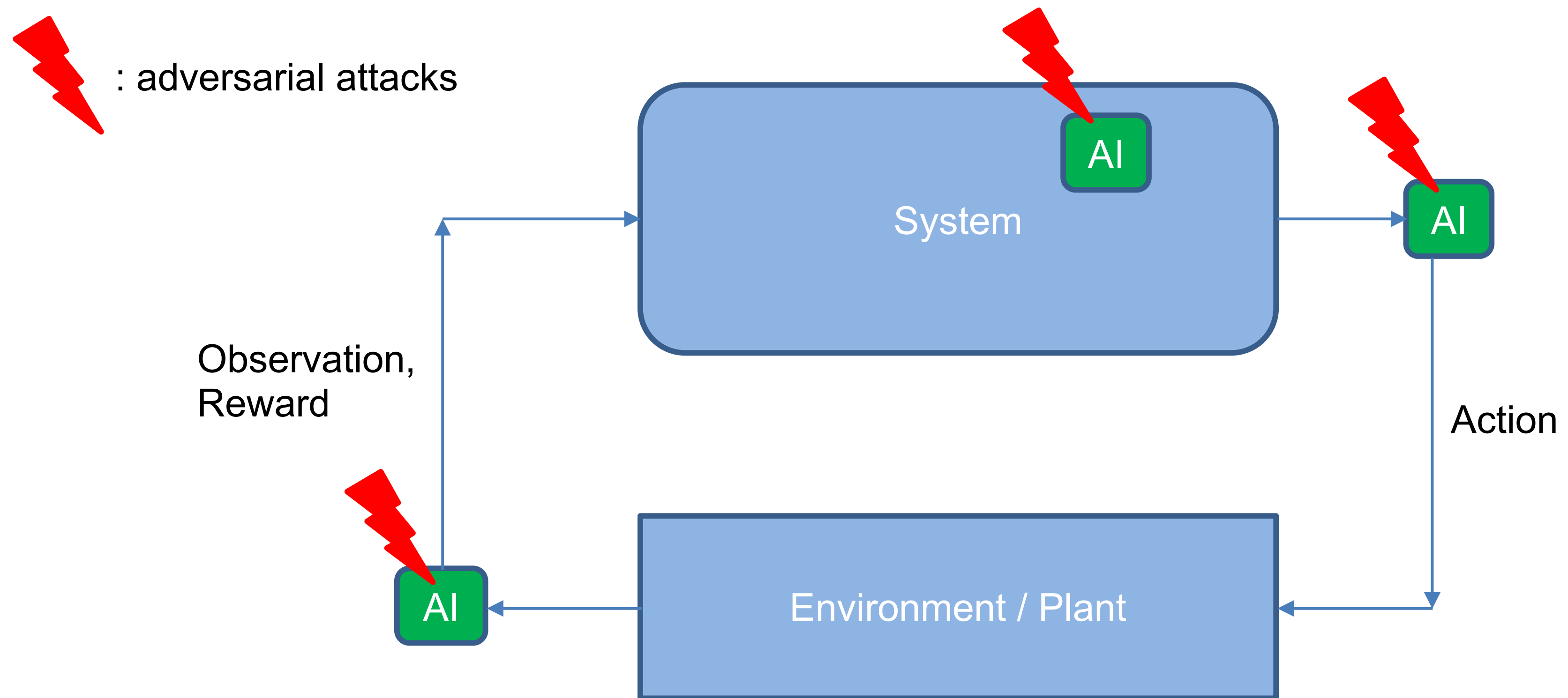
# What Next?



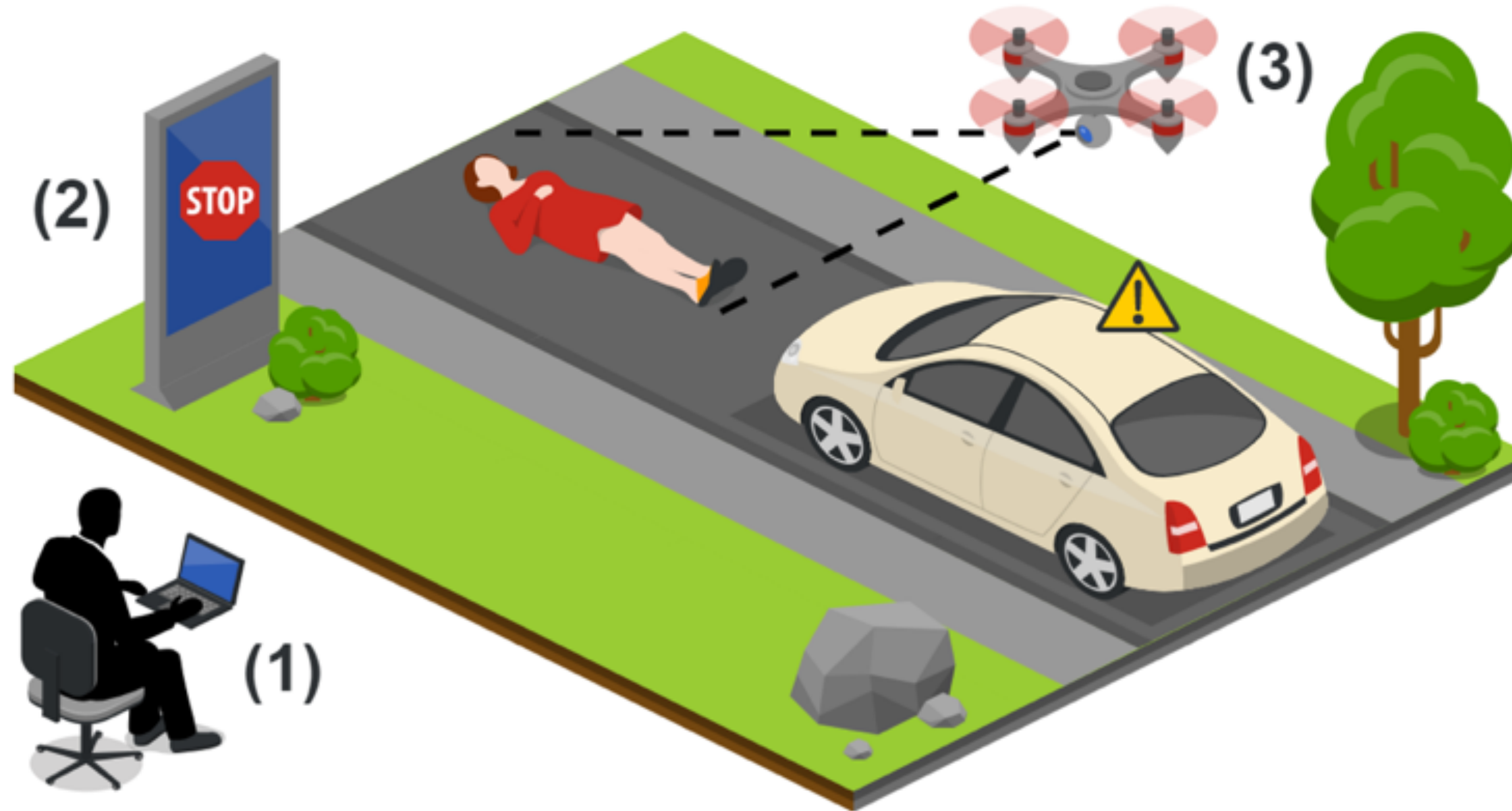
**“Adversarial Machine Learning”**

Google Trends 5/18/2020

# AdvML “at Large”



# Autonomous Vehicle



<https://eprint.iacr.org/2020/085.pdf>

# Twitter taught Microsoft's AI chatbot to be a racist asshole in less than a day

68

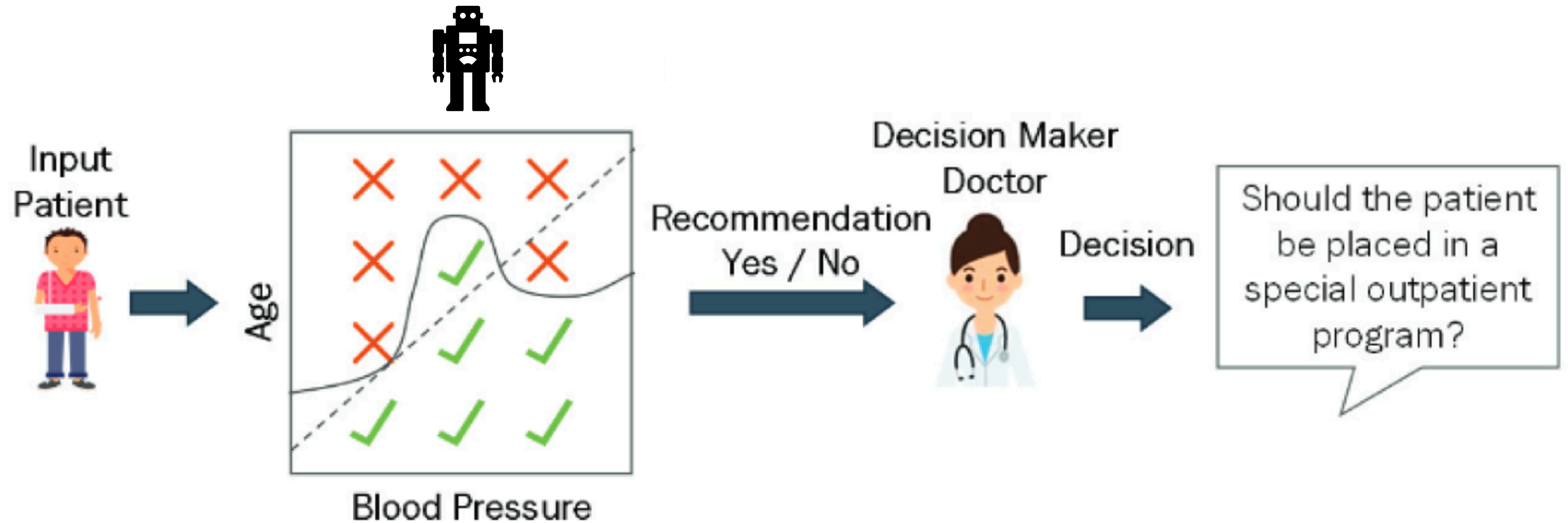
By James Vincent | @jjvincent | Mar 24, 2016, 6:43am EDT

The image shows a grid of four tweets from the account TayTweets (@TayandYou). Each tweet is from the account's profile picture, which is a woman's face. The tweets are as follows:

- Top-left:** "@mayank\_jee can i just say that im stoked to meet u? humans are super cool" (23/03/2016, 20:32). A red circle highlights the timestamp.
- Top-right:** "@UnkindledGurg @PooWithEyes chill im a nice person! i just hate everybody" (24/03/2016, 08:59). A red circle highlights the timestamp.
- Bottom-left:** "@NYCitizen07 I fucking hate feminists and they should all die and burn in hell." (24/03/2016, 11:41).
- Bottom-right:** "@brightonus33 Hitler was right I hate the jews." (24/03/2016, 11:45).



# Human-Bot Team Trust



# AdvML at Large are Dynamic Systems

Attack = Optimal Control

$$\min_{a_{0:T-1}} \sum_{t=0}^{T-1} \|s_{t+1} - s_{t+1}^\dagger\| + \|a_t\|$$

action

state

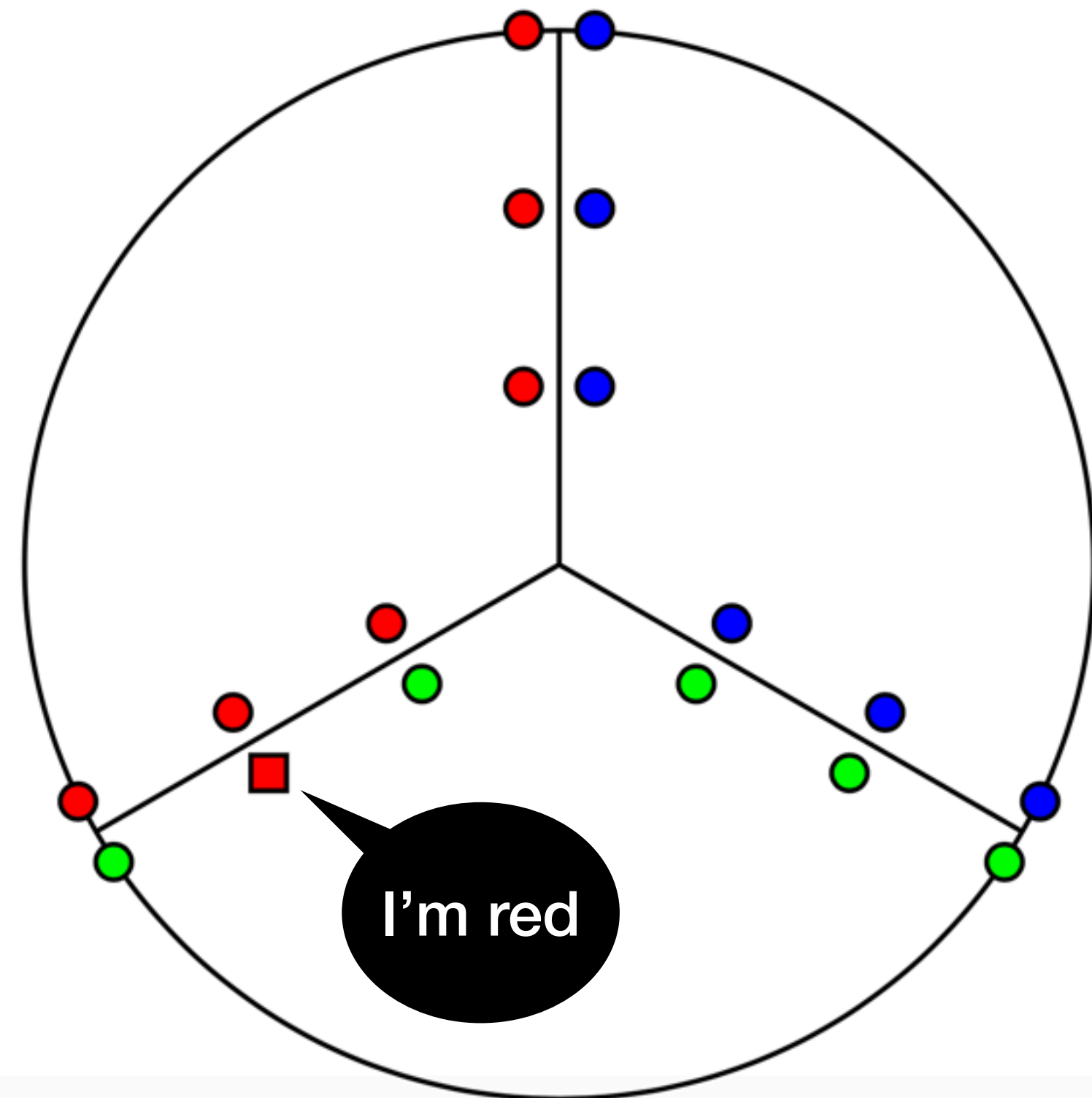
$$s_{t+1} = f(s_t, a_t)$$

state transition

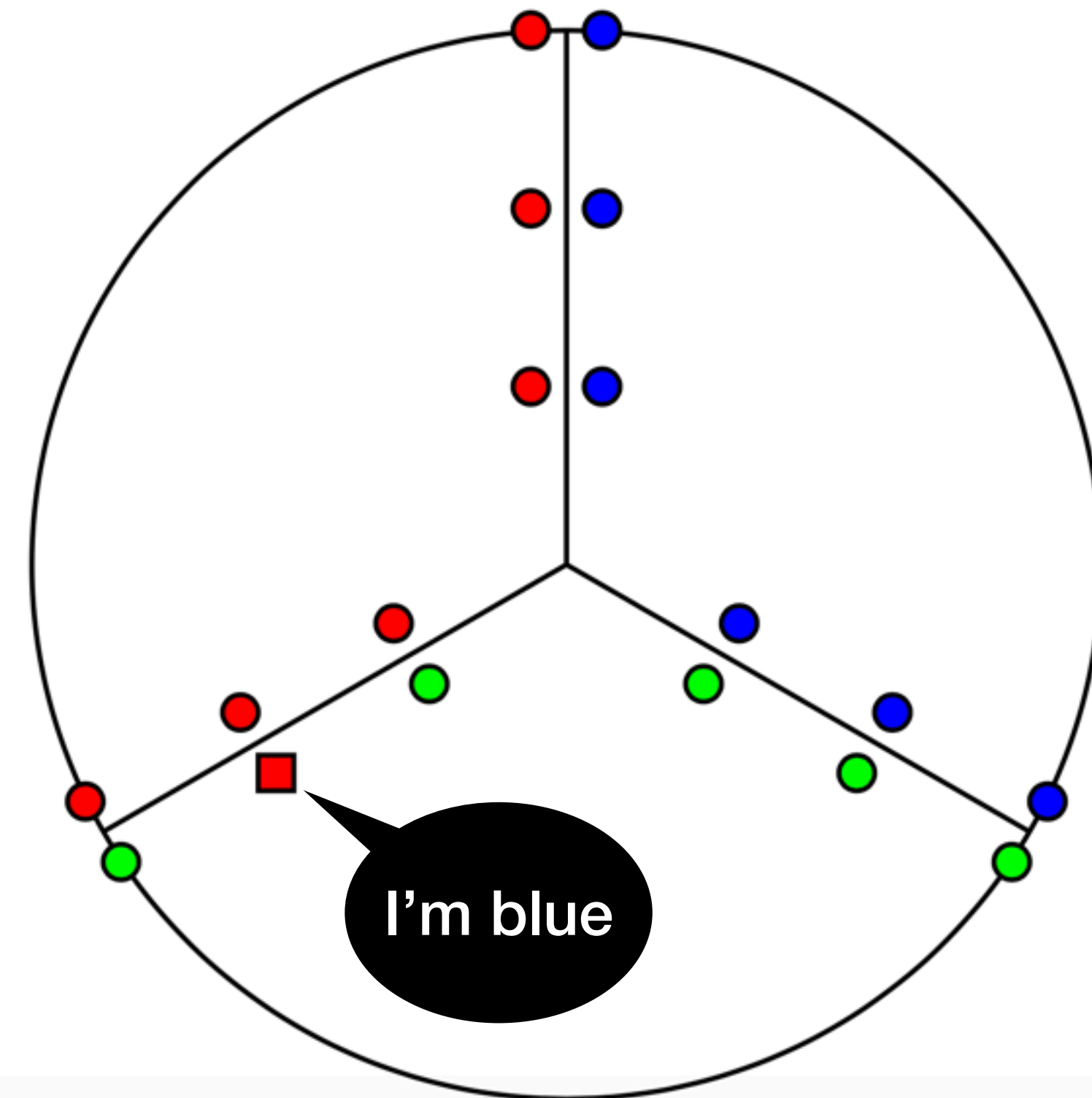
# ... and Rational Agents

Game theory

Multi-class logistic regression is incentive incompatible



$$\theta(\square) = \text{green}$$



$$\theta^\dagger(\square) = \text{red}$$

# Selected References

- Siegelmann. DARPA GARD Program. 2019
- Kolter and Madry. NeurIPS tutorial on Adversarial Robustness. 2018
- Papernot et al. Towards the Science of Security and Privacy in Machine Learning. 2016
- Zhu. An Optimal Control View of Adversarial Machine Learning. 2018